

Sequence

5910.00 101 A

PAT  
08  
3715888P

PROKARYOTIC REVERSE TRANSCRIPTASE

RELATED CASES

MW  
9-27-94

*which*  
*dit*  
This is a continuation-in-part of prior copending U.S. patent application Serial No. 07/315,427, filed February 24, 1989 and since issued as U.S. Patent No. 5,079,151 on January 7, 1992, which is a continuation-in-part of prior copending U.S. patent application Serial No. 07/315,316, filed February 24, 1989 and since issued as U.S. Patent No. 5,320,958 on June 14, 1994, which is a continuation-in-part of prior copending U.S. patent application Serial No. 07/315,432, filed on February 24, 1989 and since abandoned, which is a continuation-in-part of prior copending U.S. patent application Serial No. 07/517,946, filed on May 2, 1990, which is a continuation-in-part of prior copending U.S. patent application Serial No. 07/518,749, filed on March 2, 1990, which is a continuation-in-part of prior copending U.S. patent application Serial No. 07/753,110, filed on August 30, 1991, which is a continuation-in-part of prior copending U.S. patent application Serial No. 07/817,430, filed January 6, 1992, *31* which is a continuation-in-part of prior copending U.S. patent application Serial No. 07/979,447, filed November 20, 1992, respectively which are incorporated herein by reference.

FIELD OF THE INVENTION

The invention relates to bacterial RT enzymes which are capable of synthesizing a hybrid RNA-DNA molecule, called msDNA together with the genes which synthesize the DNA and RNA portion of the molecule.

Another aspect of the invention relates to the isolation and purification of RTs from bacterium which is capable of synthesizing msDNA. The invention deals with groups of prokaryotes

e.g., bacteria which are capable of synthesizing msDNAs by means of a reverse transcriptase. The

bacterium capable of synthesizing msDNAs is identified by testing positive by an appropriate screening test.

This is the first time that, as taught in the subject parent patent applications, reverse transcriptase has been found and isolated from a prokaryote.

5

## BACKGROUND OF THE INVENTION

Previously, there was described a chromosomal region of the bacterium Myxococcus xanthus which coded for the RNA and DNA portions of an msDNA. Dhundale et al. (Dhundale '87) "Structure of msDNA from Myxococcus xanthus: Evidence for a Long, Self-Annealing RNA precursor for the Covalently Linked, Branched RNA", Cell, Vol. 51, pages 1105-1112 (December 24, 1987). Dhundale et al. speculated that an Alu I nucleotide fragment contained all the essential coding regions to produce an msDNA. This speculation turned out to be in error.

The Alu I fragment of Dhundale et al., in fact, and inherently did not contain the gene sequence coding for an RT. The Alu I fragment was too short to code for the gene sequence coding for an RT. This was proven by way of sequence analysis by a computer program which searches for open reading frames that can potentially code for a protein. The print-out of the sequence analysis clearly shows that there is no translational reading frame in the Dhundale et al. fragment open across a stretch of DNA sufficiently long enough to encode any reverse transcriptase.

What is reported in Dhundale et al. in 1987 with respect to a bacterial reverse transcriptase was totally contrary to accepted dogma at that time about the distribution of these enzymes, i.e., that they were present only in viruses which infect eukaryotic organisms.

For the 20 years since the discovery of reverse transcriptase, it was believed that these enzymes were restricted to viruses which infect eukaryotic cells. Now, in accordance with the invention, reverse transcriptases have been identified in bacteria.

SUMMARY OF THE INVENTION

In accordance with the invention, it is shown that various bacteria have nucleotide sequences named "retrons" which encode reverse transcriptases (RTs) which are capable of synthesizing msDNAs. The invention also relates to the isolated and purified bacterial RTs. It has also been determined that the RTs of the bacteria which synthesize msDNAs possess common conserved nucleotide sequences and amino acid residues.

Representative members of the Enterobacteriaceae, Rhizobiaceae and Mycobacteriaceae families are demonstrated to be capable of synthesizing msDNA. These bacteria can be screened for the capability of synthesizing msDNA by an RT labeling or extension in vitro test.

BRIEF DESCRIPTION OF THE DRAWINGS

Figure 1 shows the restriction map of the 3.4 kb fragment around msd and downstream of msr.

Figure 2 shows the nucleotide sequence of the chromosomal region encompassing the msDNA and msd RNA coding regions and an ORF region downstream of msr and the amino acid sequence of Mx162-RT.

Figure 3 shows the amino acid sequence alignment of the msDNA-Mx162 ORF with a portion of the retroviral Pol sequences from HIV and HTLV1 and the ORF of msDNA-Ec67.

Figure 4 shows the sequence similarity of the msDNA-Mx162 reverse transcriptase with other retroelements.

Figure 5 shows the sequence comparison of the regions around the YXDD box of various reverse transcriptases.

Figure 6 shows the detection of msDNA in a clinical isolate of E. coli.

Figure 7 shows the complete primary and proposed secondary structure of msDNA-Ec67.

Figure 8 shows the determination of the RNA nucleotide sequence for the branched RNA linked to msDNA.

5 Figure 9 shows the southern blot analysis of E. coli Cl-1 Chromosomal DNA(A) and analysis of msDNA synthesis by pCl-1E and pCl-1P(B).

Figure 10 shows the restriction map of the 11.6 kb Eco RI fragment.

Figure 11 shows the nucleotide sequence of the region from the E. coli Cl-1 chromosome encompassing the msDNA, msd RNA and ORF coding regions and the amino acid  
10 sequence of Ec67-RT.

Figure 12 shows the amino acid sequence alignment of the E. coli msDNA ORF with a portion of the retroviral Pol sequence from HIV and HTLV1.

Figure 13 shows the detection of RT activity from various cell extracts.

Figure 14 shows the amino acid sequence alignment of bacterial RTs.

Figure 15 shows the nucleotide and amino acid sequence of Mx65-RT.

Figure 16 shows the nucleotide and amino acid sequence of Sa163-RT.

Figure 17 shows the nucleotide and amino acid sequence of Ec73-RT.

Figure 18 shows the nucleotide and amino acid sequence of Ec86-RT.

Figure 19 shows the nucleotide and amino acid sequence of Ec107-RT.

20 Figure 20 shows the msDNAs from total RNA prepared from each bacterial strain were specifically labeled with <sup>32</sup>P by the RT extension method (12, 14).

Figure 21 shows a collection of 63 rhizobial isolates screened for the presence of msDNA by the RT extension method.

DETAILED DESCRIPTION OF THE DRAWINGS

Figure 1. Restriction Map of the 3.4-kb fragment Around msd and Downstream of msr.

The locations and the orientation of msDNA and msdRNA are indicated by a small arrow and an open arrow, respectively. A large solid arrow represents an ORF and its orientation. The only two AluI sites (A and B) are shown and the DNA sequence between AluI (A) and AluI (B) was determined previously by Yee et al. (1984).

Figure 2. Nucleotide Sequence of the Chromosomal Region Encompassing the msDNA and msdRNA Coding Regions and an ORF Region Downstream of msr. *Seq ID NO. 1 and SEQ ID NO. 2*

The upper strand beginning at the Alu I (A) site (see Figure 1) and ending just beyond the ORF is shown. Only a part of the complementary lower strand is shown from base -301 to -600. The boxed region of the upper strand (332-408) and the boxed region of the lower strand (401-562) correspond to the sequences of msdRNA and msDNA respectively (Dhundale et al., 1987). The starting sites for DNA and RNA and the 5' to 3' orientations are indicated by open arrows. The msdRNA and msDNA regions overlap at their 3' ends by 8 bases. The circled G residue at position 351 represents the branched rG of RNA linked to the 5' end of the DNA strand in msDNA. Long solid arrows labeled a1 and a2 represent inverted repeat sequences proposed to be important in the secondary structure of the primary RNA transcript involved in the synthesis of msDNA (Dhundale et al., 1987). The ORF begins with the initiation codon at base 640. Single letter designations are given for amino acids. The YXDD amino acid sequence highly conserved among known RT proteins is boxed. Numbers on the right hand column enumerate the nucleotide bases and numbers with a\* enumerate amino acids. Small vertical arrows labeled Alu I and SmaI locate the Alu I and SmaI restriction cleavage sites, respectively. The DNA sequence was determined by the chain termination method (Sanger et al., 1977) using synthetic oligonucleotides as primer.

Figure 3. Amino acid Sequence Alignment of the msDNA-Mx162 ORF with a Portion of the Retroviral Pol Sequences from HIV and HTLV1 (and the ORF of msDNA-Ec67).

Amino acid sequences are compared with matching residues assigned as follows: (o) amino acid residues shared by all four proteins; (o) amino acid residues shared by msDNA-Mx162 and msDNA-Ec67 RTs; (x) amino acid residues shared by msDNA-Mx162 RT with HIV or HTLV1 RTs. Amino acid sequences showed are from residue-177 to -439 for HIV RT (Ratner et al., 1985); residue-15 to -277 for HTLV1 RT (Seiki et al., 1983); residue-32 to -291 for Ec-67 RT (Lampson et al., 1989); and residue-170 to -435 for Mx-162 RT (this work). The YXDD consensus sequence is outlined with a box.

Figure 4. Sequence Similarity of the msDNA-Mx162 Reverse Transcriptase with Other Retroelements.

A. Sequence similarity of the region from residue-18 to -128 of the msDNA-Mx162 RT (see Figure 2) with a carboxyl terminal region of integrase of Moloney murine leukemia virus (Mo-MLV) (residue-1070 to -1179; Shinnick et al., 1981). B. Comparison of the sequence from residue-411 to -485 of the msDNA-Mx162 RT (see Figure 2) with the sequence from residue-396 to -461 of the gap protein of human immunodeficiency virus (HIV; Ratner et al., 1985).

Figure 5. Sequence Comparison of the Regions Around the YXDD Box of Various Reverse Transcriptases.

The region from residue-304 to residue-371 of the msDNA-Mx162 RT (see Figure 2) is aligned with various RTs from different sources. The identical amino acid residues with the msDNA-Mx162 RT are indicated by open circles. The YXDD sequences are boxed. The residue numbers for the amino terminal residues and for the carboxyl terminal residues are indicated by the left and the right hand sides of the sequences, respectively. Mx-162 RT from this work (Figure 2); Ec-67 RT from Lampson et al. (1989); Ec-86 RT from Lim and Maas (1989); HIV RT from Ratner et al. (1985); HTLV1 RT from Seiki et al. (1983); Mo-MLV RT from Shinnick et al. (1981); RSV (Rous sarcoma virus) RT from Dickson et al. (1982); BLV (bovine leukemia virus) RT from Rice

Figure 6. Detection of msDNA in a clinical isolate of E. coli. Total RNA, prepared (Maniatis et al., 1982) from a 5-ml culture, was added to 50  $\mu$ l of a reaction mixture containing: 50 mM Tris-HCl (pH8.3); 6 mM  $MgCl_2$ ; 40 mM KCl; 5 mM DTT; 1  $\mu$ M dATP, dTTP and dGTP; 0.04  $\mu$ M dCTP; 0.2  $\mu$ M [ $\alpha$ - $^{32}P$ ]dCTP; and 10 units of AMV-RT (Boehringer Mannheim). The reaction mixture was incubated at 37°C for 30 min. followed by extraction with 50  $\mu$ l phenol-chloroform (1:1) and ethanol precipitation. The samples were electrophoresed on a 4% acrylamide - 8 M urea gel. Lanes: (S) molecular weight markers; MspI digest of pBR322 end-labeled with [ $\alpha$ - $^{32}P$ ]dCTP and the Klenow fragment of DNA polymerase I, (1) E. coli K-12 strain C600, (2) the same as in lane 1 except the sample was treated with RNase A (5  $\mu$ g, 10 min at 37 °C) just prior to electrophoresis, (3) clinical isolate Cl-1, (4) clinical isolate Cl-1 treated with RNase A. The clinical isolate was identified as Escherichia coli (The clinical E. coli strains were urinary tract isolates kindly provided by Dr. Melvin Weinstein from the microbiology laboratory, R.W. Johnson Hospital, New Brunswick, NJ. The clinical strain Cl-1 was identified using the API-20E identification system (API laboratory products) and gave a typical E. coli profile number of 5044552).

Figure 7. The complete primary and proposed secondary structure of msDNA-Ec67. The DNA sequence was determined by the Maxam and Gilbert method (Maxam et al., 1980) using 3'-end labeled msDNA. The RNA sequence (msdRNA; boxed region) was determined using base-specific RNases as previously described (Dhundale et al., 1987). The 2',5' Branched linkage

between the 15th rG residue and the 5' end of the DNA strand was determined using the debranching enzyme from HeLa cells as described previously (Dhundale *et al.*, 1987; Furuichi *et al.*, 1987; Ruskin *et al.*, 1985; Arenas *et al.*, 1987; the debranching enzyme was a gift from Jerard Hurwitz). The branched rG at position 15 is circled, and both RNA and DNA are numbered from their 5' ends.

Figure 8. Determination of the RNA nucleotide sequence for the branched RNA

linked to msDNA. Total RNA was prepared from the clinical strain Cl-1 and fractionated on a 5% acrylamide gel. msDNA containing full length RNA was eluted from the gel. This fraction was then labeled at the 5' end of the RNA with <sup>32</sup>P-ATP and T4 polynucleotide kinase. The 5' end labeled RNA linked to msDNA was again purified on an 18% acrylamide - 8M urea sequencing gel. The labeled RNA was then sequenced using limited digestion with base-specific RNases as described previously (Dhundale *et al.*, 1987). Lanes: OH<sup>-</sup>, partial alkaline hydrolysis ladder; (0.5 M sodium bicarbonate/carbonate pH9.2); -E, no enzyme treatment of the labeled RNA linked to msDNA; T1, RNase T1 (1U/reaction, 55°, 15 min.); U2, RNase U2 (1U and 0.5U/reaction, 55°, 15 min.); PhyM, RNase PhyM (1U/reaction, 55°, 15 min.); Bc, RNase B. cerus (2U/reaction, 55°, 15 min.); CL3, RNase CL3 (2U/reaction, 37°, 15 min.). The large gap in the sequence gel is due to msDNA linked at the rG residue at position 15 by a 2',5' phosphodiester linkage (Furuichi *et al.*, 1987). The RNA sequence at the 3'-end region from the branched rG residue (the upper part of the gel) was determined from 6% gel (data not shown).

Figure 9. Southern blot analysis of *E. coli* Cl-1 chromosomal DNA(A) and analysis of

msDNA synthesis by pLI-1E and pCl-1P(B). A: The chromosomal DNA was digested with *Eco*RI (lane 1), *Hind*III (lane 2), *Bam*HI (lane 3), *Pst*I (lane 4), and *Bgl*III (lane 5). For each lane, 3 µg of the DNA digest was applied to a 0.7% agarose gel. After electrophoresis the gel was blotted to a nitrocellulose filter, and hybridization analysis was carried out according to Southern (Southern, 1975) using msDNA labeled by AMV-RT with [ $\alpha$ -<sup>32</sup>P]dCTP as a probe. Numbers at the left represent the molecular weights in kb. B: Total DNA prepared from each strain was treated with RNase A,



separated on a 5% acrylamide gel and stained with ethidium bromide. Lane S, pBR322 digested with MspI used for molecular size markers; lane 1, DNA prepared from the host strain CL-83(recA<sup>-</sup>); lane 2, CL-83 (recA<sup>-</sup>) transformed with plasmid pCl-1E (11.6 kb EcoRI fragment; see Figure 5); lane 3, with plasmid pCl-1P (2.8-kb PstI(a)-PstI(b) fragment; see Figure 5). An arrow indicates the position of msDNA.

5

Figure 10. Restriction map of the 11.6-kb EcoRI fragment. In the Cl-1E map, the left-hand half (EcoRI to HindIII) was not mapped. In the Cl-1EP5 map, the locations and the orientations of msDNA and msdRNA are indicated by a small arrow and an open arrow, respectively. A large solid arrow represents an ORF and its orientation.

Figure 11. Nucleotide sequence of the region from the E. coli Cl-1 chromosome encompassing the msDNA and msdRNA coding regions and an ORF downstream of the msdRNA region. The entire upper strand beginning at the BalI site (see Figure 5) and ending just beyond the ORF is shown. Only a part of the complementary lower strand is shown from base 241 to 420. The long boxed region of the upper strand (249-306) corresponds to the sequence of the branched RNA (msdRNA; see Figure 7) portion of the msDNA molecule. The boxed region of the lower strand corresponds to the sequence of the DNA portion of msDNA (see Figure 7). The starting site for DNA and RNA and the 5' to 3' orientations are indicated by large open arrows. The msdRNA and msDNA regions overlap at their 3' ends by 7 bases. The circled G residue at position 263 represents the branched rG of RNA linked to the 5' end of the DNA strand in msDNA. Long solid arrows labeled a1 and a2 represent inverted repeat sequences proposed to be important in the secondary structure of the primary RNA transcript involved in the synthesis of msDNA (Dhundale et al., 1987). Note that the nucleotide at position 257 (U on the RNA transcript) and the nucleotide at position 373 (G on the RNA transcript) form a U-G pair in the stem between sequence a1 and a2. The proposed promoter elements (-10 and -35 regions) for the primary RNA transcript are also boxed. The ORF begins with the initiation codon at base 418. Single letter designations are given for amino acids. The YXDD

amino acid sequence conserved among known RT proteins is boxed. Numbers on the right hand column enumerate the nucleotide bases and numbers with a\* enumerate amino acids. Small vertical arrows labeled H and P locate the HindIII and PstI restriction cleavage sites, respectively. The DNA sequence was determined by the chain termination method (Sanger et al., 1977) using synthetic oligonucleotides as primers.

Figure 12. Amino acid sequence alignment of the E. coli msDNA ORF with a portion of the retroviral Pol sequence from HIV and HTLV1. Amino acid sequences are compared with matching residues assigned as follows: (+) amino acid common to msDNA and HIV RTs; (o) amino acid shared by msDNA and HTLV1 RTs; and (o) amino acid shared by all three proteins. Arrows divide the protein sequences into three functional domains (Toh et al., 1983; Geng et al., 1985; Varmus, 1985, Tanese et al., 1988): An amino terminal RT domain, a carboxy terminal RNase H region, and a central "tether" region. The specific amino acid residues for the RT, tether, and RNase H domains, for each protein are: HIV, 177-439, 440-600, 601-722 respectively; HTLV1, 15-277, 278-462, 463-592 respectively; msDNA ORF, 32-290, 291-465, 466-586 respectively. The YXDD polymerase consensus sequence is outlined with a box.

Figure 13. Detection of RT activity from various cell extracts. Crude cell extracts were prepared from E. coli strain C2110 (polA<sup>-</sup>) (Tanese et al., 1985; Tanese et al., 1986. E. coli strain C2110 (polA1<sup>-</sup>) was a gift from M. Roth and S. Goff) containing plasmid pCl-1EP5 encoding the msDNA-ORF (see Figure 10) as well as the vector plasmid (pUC9; Yanisch-Perron et al., 1985) alone. Extracts were also prepared from the E. coli strain PRTS7-1 (polA<sup>+</sup>) containing the cloned M-MuLV RT gene (Varmus et al., 1985; Tanese et al., 1977; Tanese et al., 1985; Tanese et al., 1986. Crude extracts were prepared essentially as described (Roth et al., 1985; Hizi et al., 1988). Crude extract equivalent to 15  $\mu$ g total protein was added to a 50  $\mu$ l reaction cocktail (50 mM tris-HCl pH7.8, 10 mM DTT, 60 mM NaCl, 0.05% NP-40, 10 mM MgCl<sub>2</sub>, 0.5  $\mu$ g poly(rC)-oligo(dG), and 0.1  $\mu$ M [ $\alpha$ -<sup>32</sup>P]dGTP and incubated at 37°C for one hour. Five  $\mu$ l of the reaction mixture was then spotted onto

5

10

Numbers on the right indicate the amino acid positions from the amino terminus for each RNA.

Sources for the sequences are Sal63 (Hsu et al. 1992b), Mx162 (Inouye et al. 1989), Mx65 (Inouye et al. 1990), Ec67 (Lampson et al. 1989b), Ec86 (Lim and Maas 1989), Ec73 (Sun et al. 1991), and Ec107 (Herzer et al. 1992).

**WEISER & ASSOCIATES**  
SUITE 500  
230 SO. FIFTEENTH ST.  
PHILADELPHIA, PA 19102  
(215) 875-8383  
TELECOPIER (215) 875-8394

shown from bases 121-300. The boxed region of the upper strand (positions 143-191) and the boxed region of the lower strand (positions 186-250) correspond to the sequences of msdRNA and msDNA, respectively. The starting sites for DNA and RNA and the 5' to 3' orientation are indicated by open arrows. The msdRNA and msDNA regions overlap at their 3' ends by 6 bases. The circled G residue at position 206 represents the branched guanosine of RNA linked to the 5' end of the DNA strand in msDNA. Long solid arrows labeled a1 and a2 represent inverted repeat sequences proposed to be important in the secondary structure of the primary RNA transcript involved in the synthesis of msDNA. The ORF begins with the initiation codon at base 279. The YXDD amino acid sequence highly conserved among known RT proteins is boxed. Numbers on the right-hand column enumerate the nucleotide bases, and numbers with asterisks enumerate amino acids (single-letter code). The DNA sequence was determined by the chain-termination method using synthetic oligonucleotides as primers.

Figure 16 shows nucleotide sequences of 3,060 bases encompassing msr, msd, and the RT gene of S. aurantiaca. The sequence from base 421 to base 720 which contains msr and msd is shown double stranded. The boxed regions of the upper strand (bases 440 to 540) and the lower strand (bases 508 to 670) correspond to the sequences of msdRNA and msDNA, respectively. The starting sites for msDNA and msdRNA are indicated by open arrows. The circled G at the position 458 is the branched rG of msdRNA linked to the 5' end of msDNA. Long solid arrows labeled with a1 and a2 represent inverted repeated sequences proposed to form the secondary structure in the primary RNA transcript which serves to prime msDNA synthesis. Amino acids are indicated by single letters. The YXDD sequence highly conserved among known RTs is boxed. X<sup>e</sup> and B<sup>f</sup> sites are indicated by arrows. Numbers on the right-hand side and numbers with asterisks represent numbers for bases and amino acids, respectively.

Figure 17 shows the sequences of msdRNA and msDNA which are boxed and their orientations are indicated by open arrows. The branched G residue at position 10425 is circled. The

inverted repeat sequences require for the biosynthesis of msDNA - Ec73 are shown by arrows labeled a1 and a2. Amino acid residues of Ec73-RT are shown by a single-letter code put at the center of each codon.

Figure 18 shows the restriction map of the 3.5 kb insert of pDB808 and nucleotide sequence of chromosomal determinants of the msDNA-RNA compound of E. coli B. (A) Restriction map of the 3.5 kb insert of clone pDB808. The solid bar represents the region whose sequence is presented in (B). Transcription is from left to right. Restriction enzymes are: P, PstI, H<sub>u</sub>, HpaI; B, BglII; X, XhoI. (B) Nucleotide sequences of the chromosomal determinants. Only the strand corresponding to the transcript is shown. Nucleotides are numbered starting from the first base observed in the msdRNA. The msdRNA coding region is overlined, and the msDNA coding region is underlined. The msDNA sequence is complementary to the sequence shown in this figure. Inverted repeats are indicated by double-dashed lines. The G at position 14 is the branched guanylate of msdRNA in the msDNA-RNA compound. IR, 12 bp inverted repeat.

Figure 19 shows sequence of the retron and flanking regions of Ec107. The sequences corresponding to the K-12 genomic DNA are shown in lower case letters from bases 1-99 and 1400-1540. The msRNA and msDNA regions are boxed. Also indicated are the a1-a2 conserved inverted repeats (indicated by arrows) and the branched G, which is circled. The RT consists of 319 amino acids and contains the YXDD sequence (boxed) which is highly conserved among known RTs. The transcription start site occurs at base 170; a possible terminator is indicated by head-to-head arrows following the RT coding region. Primer extension was utilized in order to determine the transcription start site. These sequence data will appear in the EMBL/GenBank/DDJB Nucleotide Sequence Data Libraries under the accession number X62583.

DETAILED DESCRIPTION OF THE INVENTION

The description which follows describes msDNA and RT from Myxococcus xanthus. This is a typical bacterium which belongs to a genus of bacteria, whose representative members possess an RT capable of synthesizing msDNA.

5 The existence of a peculiar branched RNA-linked DNA molecule called msDNA (multicopy single-stranded) has been demonstrated in various myxobacteria, Gram-negative soil bacteria (Yee et al., 1984; Dhundale et al., 1985; Furuichi et al., 1987a,b; Dhundale et al., 1987; Dhundale et al., 1988b). msDNA (msDNA-Mx162) from Myxococcus xanthus consists of 162-base single stranded DNA, the 5' end of which is linked to the 2' position of the 20th rG residue of a 77-  
10 base RNA molecule (msdRNA) by a 2', 5'-phosphodiester linkage (Dhundale et al., 1987). It exists at a level of approximately 700 copies per genome. Stigmatella aurantiaca also possesses an msDNA (msDNA-Sal63) which is highly homologous to msDNA-Mx162 (Furuichi et al., 1987b). In addition to msDNA-Mx162, M. xanthus has another smaller species of msDNA (mrDNA or msDNA-Mx65), which has no primary sequence homology with msDNA-Mx162 or msDNA-Sal63 (Dhundale et al., 1988b). However, all msDNAs so far characterized share key structural features such as a branched  
15 rG residue, stem-and-loop structures in RNA and DNA molecules, and a DNA-RNA hybrid at the 3' ends of DNA and RNA molecules.

Previously it was predicted that reverse transcriptase is required for msDNA biosynthesis on the basis of the finding that msdRNA is derived from a much longer precursor, which  
20 can form a very stable stem-and-loop structure (Dhundale et al., 1987). This precursor molecule was proposed to serve as a primer for initiating msDNA synthesis as well as a template to form the branched RNA-linked msDNA. The latter reaction requires reverse transcriptase activity. In M. xanthus, the region coding for the RNA molecule (msr) is located on the chromosome in the opposite orientation to the msDNA coding region (msd) with the 3' ends overlapping by 6 bases for msDNA-Mx65 (Dhundale et al., 1988b) or by 8 bases for msDNA-Mx162 (Dhundale et al., 1987). In addition,  
25 as in all the msDNAs found in myxobacteria, there is an inverted repeat comprised of a 14-base

sequence for msDNA-Mx65 (Dhundale et al., 1988b) or a 34-base sequence for msDNA-Mx162 (Dhundale et al., 1987) and a 33-base sequence for msDNA-Sal63 (Furuichi et al., 1987b) immediately upstream of the branched G residue and a sequence immediately upstream of the msDNA coding region. As a result of this inverted repeat, a longer primary transcript beginning upstream of the RNA coding region and extending through the msDNA coding region is considered to self-anneal and form a stable secondary structure. When three base mismatches were introduced into the secondary structure immediately upstream of the branched rG residue, msDNA synthesis was almost completely blocked. However, if three additional base substitutions were made on the other strand to resume the complementary base pairing, msDNA production was restored (Hsu et al., 1989). This result strongly supports the proposed model for msDNA synthesis.

It was also shown that a deletion mutation at the region 100 base pairs (bp) upstream of the DNA coding region (msd) and an insertion mutation at a site 500 bp upstream of msd caused a significant reduction in msDNA production (Dhundale et al., 1988a). This indicates that there is a cis- or trans-acting positive element required for msDNA synthesis in this region. In this report we determined the DNA sequence of this region and found an opening reading frame (ORF) coding for 485 amino acid residues beginning with an initiation codon, ATG, which is located 77 bp upstream of msd (or 231 bp downstream of msr). The very close proximity between msd and the ORF suggests that they may be transcribed as a single transcript. The amino acid sequence of the ORF shows similarity with retroviral reverse transcriptases. We discuss a possible origin of the reverse transcriptase gene as well as a possible relationship between the msDNA system and retroviruses. Recently, some strains of Escherichia coli were found to produce msDNA and the gene for reverse transcriptase which is essential for msDNA production, is linked to the msd region, (Lim and Maas, 1989; Lampson et al., 1989b). Comparison of the msDNA systems of M. xanthus and E. coli raises an intriguing question as to how the extensive diversity found in msDNA systems has emerged in bacteria and what possible functions msDNA may have.

In a preceding paper, it was demonstrated that msDNA is in fact synthesized by reverse transcriptase in a cell-free system in M. xanthus (Lampson et al., 1989a).

Reverse transcriptases are isolated, and if desired, purified, and biological characterization carried out, if desired, by known methods such as those described in Lampson, B.C., M. Viswanathan, M. Inouye and S. Inouye, "Reverse Transcriptase from *Escherichia coli* Exists as a Complex with msDNA and is Able to Synthesize Double-stranded DNA", *J. Biol. Chem.* 265: 8490-8496 (1990), which is incorporated by reference as if fully set forth herein.

## RESULTS AND DISCUSSION

### Identification of an ORF Associated with msd

On the basis of mutations closely associated with msd which significantly reduce msDNA production, it was assumed that in this region there is a cis- or trans-acting element which is essential for msDNA synthesis (Dhundale *et al.*, 1988a). Figure 1 shows a restriction map around msd. The msDNA coding region is shown by a thin arrow from right to left (msd), and the msdRNA coding region by a thick open arrow (msr). In the previous work (Dhundale *et al.*, 1988a), two mutations were constructed; one, a deletion mutation in which the sequence from Alu I(b) to SmaI was replaced by a gene for kanamycin resistance (see Figure 1), and the other an insertion mutation at the SmaI site by a gene for kanamycin resistance (see Figure 1).

In order to elucidate the properties of the element required for msDNA production, the DNA sequence of the region upstream of msd was determined as shown in Figure 2. A long open reading frame (ORF) beginning with an initiation codon was found 77 bases upstream of msd. The ORF is preceded by a ribosome binding sequence of AGG (residue 630 to 632) 7 bases upstream of the initiation codon. The ORF codes for a polypeptide of 485 amino acid residues. The Alu I(b) and SmaI sites (see Figure 1), where mutations inhibiting msDNA synthesis were created, are located at amino acid residue -12 and -142 of the ORF, respectively or at the nucleotide sequence from residue -672 to -675, and from residue -1061 to -1066, respectively (Figure 2). In Figure 2, msd or the DNA sequence corresponding to the msDNA sequence is indicated by the closed box on the lower strand and the orientation is from right to left. Similarly, the msdRNA sequence (msr) is also indicated by



the closed box on the upper strand and the orientation is from left to right. The msd and msr regions overlap by 8 bases. An inverted repeat is also indicated by arrows with letters a1 and a2. This inverted repeat comprises a 34-base sequence immediately upstream of the branched G residue (residue 317 to 350; sequence a2 in Figure 2) and another 34-base sequence at the 3' end (residue 597 to 564; sequence a1). This inverted repeat is essential to form a stem structure which provides a stable secondary structure in a long primary transcript. This secondary structure is considered to serve as the primer as well as the template for msDNA synthesis (Dhundale et al., 1987; Hsu et al., 1989).

#### Sequence Similarity with Retroviral Reverse Transcriptases

When the amino acid sequence of the ORF was compared with known proteins, a striking similarity was found between the sequence from Leu-308 to Ser-351 and retroviral reverse transcriptases (RT). In particular, this region contains the YXDD sequence, the highly conserved sequence in all known RTs. This sequence (Tyr-344 to Asp-347) is boxed in Figure 2. In Figure 3, the ORF sequence of 266 amino acid residues from Ala-170 to Lys-435 is compared with RTs from HIV (human immunodeficiency virus; Ratner et al., 1986) and HTLV1 (human T-cell leukemia virus type 1; Seiki et al., 1983). As mentioned above, within the sequence of 44 amino residues from Leu-308 to Ser-351, there are 14 and 12 identical residues with HIV (32%) and HTLV1 (27%), respectively. The entire RT domains of HIV and HTLV can also be aligned with the ORF sequence from Ala-170 to Lys-435, with much less similarity as shown in Figure 3. However, the same region was found to be extremely well aligned with the RT which was recently found in a clinical strain of Escherichia coli (Lampson et al., 1989b). This E. coli RT consists of 586 amino acid residues, and its amino terminal domain (residue-32 to -291) and the carboxyl terminal domain (residue-466 and -586) have been demonstrated to have sequence similarity with retroviral RT and ribonuclease H. This RT gene from E. coli was shown to be required for the production of msDNA (msDNA-Ec67) and to have reverse transcriptase activity (Lampson et al., 1989b). Figure 3 shows that the sequence similarity between E. coli and M. xanthus RTs is distributed within almost the entire RT region; in particular in the region from Tyr-181 to Ser-212, 15 out of 32 residues are identical (47% similarity);

in the region from Gly-226 to Gly-265, 19 out of 40 residues (48% similarity); in the region from Leu-308 to Ser-351, 26 out of 44 residues (59% similarity); and in the region from Lys-354 to Asn-408, 21 out of 55 residues (38% similarity). Overall, similarity from Ala-170 to Lys-435 is 32% (85 out of 266 residues are identical). In spite of these similarities, the M. xanthus ORF does not have the domain, which shows apparent sequence similarity with ribonuclease H (RNase H). The RNase H domain is found to be located in the carboxyl terminal region of the same polypeptide in which the RT domain exists in the amino terminal region in the case of the E. coli RT and other retroviral RTs. In the preceding paper, it was shown that there is a precise coupling between RT and RNase H activity (Lampson et al., 1989a). Therefore, RNase H may still reside with the ORF, or RNase H may be encoded by a separate gene.

#### Sequence Similarity with Other Proteins

In contrast to the E. coli RT and other retroviral RTs, the ORF found in M. xanthus has a long amino terminal extra domain consisting of approximately 170 residues. Interestingly, this region shows some sequence similarities with the carboxyl terminal region associated with integration protein of Mo-MLV (Moloney murine leukemia virus; Shinnick et al., 1981) (see Figure 4A); the sequence from Pro-18 to Leu-128 of the ORF shows 22% similarity (24 out of 111 residues) with the region from Pro-1070 to Leu-1179 of the gag-pol polyprotein of Mo-MLV. It should be noted that this region of Mo-MLV is unique for Mo-MLV integration protein and does not share sequence similarity with other retroviral endonucleases (Johnson et al., 1986). It is also interesting to notice that in Ty retrotransposon, this domain is located in front of the RT domain in contrast to the retroviral endonuclease domain (Clare and Farabaugh, 1985).

As pointed out above, the ORF does not have homology to E. coli or retroviral RNase H. Instead, it has a short sequence of approximately 80 residues after the RT domain. In this region, one can also find sequence similarity with a part of the gag region of HIV. As shown in Figure 4B, the sequence from Gly-411 to Glu-485 has 22 identical amino acid residues (31% similarity) with the region from Gly-396 to Pro-461 of the gag protein of HIV (Ratner et al., 1985).

Requirement of Reverse Transcriptase

The fact that disruption of the ORF significantly reduced msDNA production in M. xanthus (Dhundale et al., 1988a) and the fact that the ORF has sequence similarity with retroviral RTs strongly supports the previous hypothesis that RT is required for the synthesis of msDNA (Dhundale et al., 1987). Recently, we were able to demonstrate that msDNA is indeed synthesized by reverse transcriptase activity in a cell-free system (Lampson et al., 1989a). The fact that a small amount of msDNA (3% of the wild type level) is still produced in the ORF mutants (Dhundale et al., 1988a) is most likely due to another RT associated with smaller msDNA (msDNA-Mx65; previously assigned mrDNA; Dhundale et al., 1988b). In fact, an ORF has been found to be associated with the region responsible for msDNA-Mx65 production.

At present it is unknown if the ORF is transcribed together with msdRNA from a common upstream promoter or if the ORF has its own independent promoter. Previously, a major RNA transcript of approximately 375 bases by S1 mapping (Dhundale et al., 1987) was identified. This transcript covers the region from approximately 75 bases upstream of msr (at around residue-256 in Figure 2) to approximately 70 bases upstream of msd (at around residue-632 in Figure 2). This indicates that this RNA transcript ends at the ribosome binding site (AGG, 630-632) of the ORF. It is possible that the primary RNA transcript covers not only the msr-msd region but also the entire ORF. This transcript of approximately at least 2 kilobases (kb) is then used as the mRNA for the ORF to produce RT. At the same time, the 5' untranslated region of 350 bases forms a stable secondary structure which serves as a primer and a template for msDNA synthesis as previously proposed (Dhundale et al., 1987). Because of the secondary structure, the 5' end region is probably much more stable than the ORF mRNA region. As a result, only the 375-base RNA from the 5' end of the transcript was detected in the previous work. In E. coli, the RT gene was shown to be transcribed from a single promoter for the msr region (Lampson et al., 1989b).

Evolution of Reverse Transcriptase

All of the RTs so far identified are from eukaryotic origins, and associated with either retroviruses or retrotransposons. DNA synthesis for retroviruses and transposition events for retrotransposons occur via RNA which is used as a template for RTs (see review by Varmus, 1985).

5 From amino acid similarity in various RTs, possible evolutionary relationships among these RTs has been proposed (Yuki et al., 1986).

The present invention demonstrates that RTs are not specific to eukaryotes but exist in prokaryotes as well. An intriguing question arises as to the evolutionary relationship between prokaryotic and eukaryotic RTs and the origin of RT. In order to compare the amino acid sequences of these RTs, the sequence of the M. xanthus RT from Gly-304 to Leu-371 was chosen, since this sequence includes the YXDD box, the most conserved region among different RTs. In Figure 5A this sequence is compared with 13 other representative RTs from bacteria, yeast, plant, mitochondrial plasmid, and animal retroviruses. Within these 14 sequences, the D-D sequence (residues-346 and -347) is completely conserved, and both G-311 and Y-344 are also well conserved except for Ty-RT. Besides these residues, L-308, P-309, Q-310, S-315, P-316, L-330, S-351, and L-371 are fairly well conserved among these sequences. On the basis of the numbers of identical amino acid residues, M. xanthus RT has the following similarities with other RTs: 47% (32 amino acid residues) with E. coli C1-1 RT; 41% (28) with E. coli B RT; 24% (16) with HIV, BLV, and mitochondrial plasmid RTs; 22% (15) with Mo-MLV RT; 21% (14) with RSV, 17.6, gypsy, and Tal-3 RTs; 19% (13) with HTLV1 RT; 15% (10) with Ty912 RT; and 9% (6) with Copia RT. On the basis of the phylogenetic relationships among RTs proposed by Yuki et al. (1986), and the present data, a dendrogram of homology of various RTs may be constructed as shown in Figure 5B. As proposed earlier (Yuki et al., 1986), modern RTs are composed to two major groups I and II. One group (group II) consists of retrotransposons found in yeast (Ty912), plant (Tal-3), and Drosophila (Copia). Bacterial RTs seem to belong to the other group (group I) together with other retrotransposons from Drosophila such as 17.6 and gypsy, mitochondrial plasmid RT, and retroviral RTs. This indicates that both prokaryotic and eukaryotic RT genes were possibly derived from a single ancestral RT gene.

Origin of the *M. xanthus* Reverse Transcriptase

In addition to the sequence similarity between the *M. xanthus* RT and RTs from retroviruses and retrotransposons, msDNA shares other interesting similarities with retroviruses and retrotransposons; msDNA (synthesis of single-stranded DNA) starts at a site 77 bases upstream of the RT gene and the orientation of DNA synthesis is opposite to the direction of translation of the RT gene. In the case of retroviruses and retrotransposons, single-stranded DNA synthesis proceeds at the 5'-end untranslated region of an RNA molecule which serves as the mRNA for RT as well (Weiss *et al.*, 1985). The orientation of DNA synthesis is also opposite to the direction of translation of the RT gene. In the case of msDNA synthesis an RNA transcript itself serving as a template also serves as a primer by self-annealing to form a stable secondary structure (Dhundale *et al.*, 1987), whereas in the case of retroviruses and retrotransposons tRNAs are recruited from the cell for the priming reaction. At present it is unknown if branched RNA-linked msDNA is the final product of an unknown function or if it is a stable intermediate leading to other products.

Furthermore, it is of great interest whether the *M. xanthus* RT is associated with a complex such as virus-like particles such as those found for yeast Ty1 element (Eichinger and Boeke, 1988). In a preliminary experiment, msDNA of *M. xanthus* exists as a complex with proteins in the cell which sediments as a 22S particle. Characterization of this complex may shed light on questions concerning the relationship between msDNA and retrocomponents as well as the functions of msDNA.

At present, there is no information to support the possibility that msDNA may be a transposable element or an element associated with a provirus (or prophages). It is important to point out that the RT gene from *M. xanthus* appears to be as old as other genomic genes for the following reasons: (a) Nine independent natural isolates of *M. xanthus* from various sites (including Fiji Island and eight different sites in the United States) contained mutually hybridizable msDNA (Dhundale *et al.*, 1985). Since under the same hybridization condition, msDNA-Mx162 did not hybridize with msDNA-Sa163 [which has extensive homology in both DNA and RNA sequences with msDNA-Mx162; Dhundale *et al.*, (1987)], the nine independent strains *M. xanthus* are assumed to contain almost identical msDNA. (b) The codon usage of the Mx-162 RT is almost identical to those found

in other M. xanthus genes (Table 1). M. xanthus is known to have a very high G+C content (70%; Johnson and Ordal, 1968) and as a result, all the genes so far characterized have very high G+C contents at the third positions of codons used; 85.4% for vegA (Komano et al., 1987), 85.7% of ops (Inouye et al., 1983), 87.2% for tps (Inouye et al., 1983), 88.4% for mbhA (Romeo et al., 1986), and 93.9% for sigma factor. The average G+C content of the third positions is calculated to be 90.0% for these genes (Table 1). Surprisingly, the G+C content of the third positions of the RT codons is highest among these genes (95.5%; Table 1).

In contrast, the E. coli msDNA system including the RT gene is considered to have been acquired much later in the evolution of E. coli. Reasons for this conclusion include: (a) Only four strains out of 89 independent clinical E. coli strains were found to produce msDNAs (Lampson et al., 1989b). (b) The codon usage of the E. coli RT is significantly different from the general codon usage of E. coli genes obtained from 199 E. coli genes (Maruyama et al., 1986). In particular, out of 62 arginine codons used in the E. coli RT, 40 (65%) use AGA or AGG in contrast to 2.7% for the AGA+AGG usage among all arginine codons in 199 E. coli genes (see Table 1). The AGA and AGG codons are the least used codons in E. coli (Maruyama et al., 1986). In addition to AGA and AGG codons, many other codons, GCC and GCG for Ala, CGU and CGC for Arg, CAG for Gln, GGC and GGA for Gly, CAC for His, AUC and AUA for Ile, UUA, CUU and CUG for Leu, UUC for Phe, CCU and CCG for Pro, UCG for Ser, ACC and ACA for Thr, and GUC for Val. (c) Although the E. coli msDNAs share little sequence homology, they all share the key secondary structures of a branched rG residue, a DNA-RNA hybrid at the 3' ends of the msDNA and msdRNA, and stem-and-loop structures in RNA and DNA strands (Lampson et al., 1989b; Lim and Maas, 1989).

These results clearly demonstrate distinct differences between the msDNA systems of E. coli and M. xanthus. Myxobacteria are common organisms in soil and are found all over the world regardless of climate, and considered to diverge from their nearest bacterial relatives about  $2 \times 10^9$  years ago when the atmosphere became aerobic (see a review by Kaiser, 1986). Since it is reasonable to assume that the M. xanthus RT gene is as old as other genomic genes, the RT gene existed much

before eukaryotic cells appeared ( $1.5-0.9 \times 10^9$  years ago). The relatedness between various

prokaryotic and eukaryotic RTs as shown in Figures 5A and B strongly supports the existence of a single ancestral gene for all RTs. It is possible that such an ancestral RT gene was independently recruited into different systems such as the msDNA system, the retrotransposon system, and the retroviral system. Alternatively, the msDNA system may be a primitive ancestral system from which retrotransposons and retroviruses originated. In this regard, it is intriguing to point out other sequence similarities between the M. xanthus RT-ORF and other retroelements (see Figure 4) other than RT itself as well as the similar mode of initiation of DNA synthesis by RT as discussed earlier.

At present, it is beyond our speculation why the E. coli msDNA systems are so diverged in contrast to the M. xanthus msDNA system and how they were acquired into the genomes of some E. coli strains. However, it should be noted that the E. coli RTs are most related to the M. xanthus RT indicating that they were not derived from eukaryotic origins. Possible origins of retroviruses have been discussed (Temin, 1980). The recent finding of an imposon in a genetic component for a mouse gene also raises an interesting question concerning the evolution of retroelements (Stavenhagen and Robins, 1988). Further characterization of the prokaryotic RTs and the msDNA systems will provide clues to the origins of RT and other retroelements.

## EXPERIMENTAL PROCEDURE

### DNA Manipulation and Plasmids

DNA manipulation was performed as described by Maniatis et al. (1982). The plasmid isolation was as originally described by Birnboim and Dolly (1979). Plasmid pmsSB7 containing the 5 kb Sall-BamHI fragment shown between the Sall and BamHI sites of pUC9 (Vieira and Messing, 1982) was used. After the 2.2 kb Sall-SmaI fragment from pmsSB7 was subcloned between the Sall and SmaI sites of pUC9, all RsaI fragments were gel-purified and cloned into pUC9 for DNA sequence.

DNA sequence

DNA sequence was determined by the chain termination method (Sanger et al., 1977) using single-stranded or double-stranded DNA as templates with synthetic oligonucleotides.

Other Material and Methods

5 Restriction enzymes were purchased from either Bethesda Research Laboratories or New England BioLabs. [ $\alpha$ -<sup>35</sup>S] dATP was from Amersham. Sequenase, Version 2.0 Kit was purchased from United States Biochemical Corporation for DNA sequences.

Cyborg program from International Biotechnologies Inc. was used to search sequence homology in GenBank Release 55.

10 Screening of bacteria for retron synthesized msDNAs was performed by the methods of Lampson et al. J. Bacteriol., 173:5363-5370 (1991), or Yee et al., Cell, 38, 203-209 (1984).

RTs were identified and isolated by the method of Lampson et al., J. Biol. Chem., 265:8490-8496.

msDNA in Escherichia coli

15 The recent serendipitous finding of msDNA (msDNA-Ec86) in E. coli B by Dongbin Lim and Werner Maas (D. Lim et al., 1989) prompted a to search for msDNA in other E. coli strains. Previously established by Yee et al. (T. Yee et al., 1984), msDNA is not found in the common laboratory strain K12, however, to our surprise, it was in a clinical E. coli strain isolated from a patient with a urinary tract infection. Fifty independent E. coli urinary tract isolates were examined

20 for the presence of msDNA (The clinical E. coli strains were urinary tract isolates kindly provided by Dr. Melvin Weinstein from the microbiology laboratory, R.W. Johnson Hospital, New Brunswick, NJ. The clinical strain Cl-1 was identified using the API-20E identification system (API laboratory products) and gave a typical E. coli profile number of 5044552.). The screening method involved treatment of total RNA prepared from each strain with (AMV) RT in the presence of [ $\alpha$ -<sup>32</sup>P]dCTP

plus dATP, dTTP, and dGTP followed by polyacrylamide gel electrophoresis. Since msDNA contains



a DNA-RNA duplex structure, the 3' end of the DNA molecule serves as an intramolecular primer and the RNA molecule as a template for RT. When RNA prepared from one of the clinical strains, E. coli Cl-1, was labeled in this manner, two distinct, low molecular weight bands of about 160 bases became labeled with  $^{32}\text{P}$  and are shown in Figure 6. If the labeled sample is digested with ribonuclease (RNase) A prior to loading on the gel, a single band corresponding to 105 bases of single-stranded DNA is detected (lane 4). This indicates that both bands in lane 3 contain a single-stranded DNA of identical size. The two labeled bands observed prior to RNase treatment (lane 3) are due to two species of msDNA comprised of a single species of single-stranded DNA linked to RNA molecules of two different sizes. RNA molecules of two different sizes have been observed at the 5' ends of msDNA from myxobacteria in which a precursor molecule contains a longer RNA which is processed into a smaller mature form (Dhundale et al., 1987; Furuichi et al., 1987). Among the 89 clinical isolates screened, three other strains produced msDNA-like molecules of varying size and quantity, suggesting extensive diversity among these molecules. As previously reported (Dhundale, 1985), msDNA was not observed in the E. coli K-12 strain, C600 (lanes 1 and 2, Figure 6).

#### Nucleotide sequence of msDNA Ec-67

To determine the base sequence of the DNA molecule, the RNA-DNA complex isolated from the clinical strain was labeled at the 3' end of the DNA molecule with AMV-RT and  $[\alpha\text{-}^{32}\text{P}]\text{dATP}$ . By adding dideoxy-CTP, ddTTP, and ddGTP to the reaction mixture, a single labeled adenine is added to the 3' end of the DNA molecule. RNA is removed with RNase A+ T1 and the end-labeled DNA is subjected to the Maxam and Gilbert sequencing method (Maxam et al., 1980). Figure 7 shows that msDNA consists of a single-stranded DNA of 67 bases and, as in the case of msDNAs from myxobacteria (Yee, 1984; Dhundale, 1987), it can form a secondary hair-pin structure. The primary sequence, however, is not homologous to any of the myxobacterial msDNAs, nor to the msDNA from E. coli B (msDNA-Ec86; Lim and Maas, personal communication).

The sequence of the RNA molecule was determined using the RNA-DNA complex purified from *E. coli* Cl-1. The RNA sequence was determined using base specific RNases as described previously (Dhundale *et al.*, 1988). As shown in Figure 8, a large gap is observed in the RNA sequence "ladder". This gap is due to the DNA strand branched at the 2' position of the 15th rG residue of the RNA strand which produces a shift in mobility of the sequence ladder (see Figure 7). The RNA consists of 58 bases with the DNA molecule branched at the G residue at position 15 by a 2',5'-phosphodiester linkage. The branched G structure was determined as previously described for msDNAs from myxobacteria (Dhundale, 1987; Furuichi *et al.*, 1987). After RNase (A and T1) treatment, msDNA retains a small oligoribonucleotide linked to the 5' end of the DNA molecule due to the inability of RNases to cleave in the vicinity of the branched linkage. The 5' end was labeled with [ $\gamma$ - $^{32}$ P]ATP using T4 polynucleotide kinase and the labeled RNA molecule was detached from the DNA strand by a debranching enzyme purified from HeLa cells (Ruskin *et al.* 1985; Arenas *et al.*, 1987; the debranching enzyme was a gift from Jerard Hurwitz). This small RNA was found to be a tetranucleotide which could be digested with RNase T1 to yield a labeled dinucleotide (not shown). Since RNase T1 could not cleave the RNA molecule at the G residue before debranching enzyme treatment, it was concluded that the single-stranded DNA is branched at the G residue via a 2',5'-phosphodiester linkage. In addition, partial RNase U<sub>2</sub> digestion cleaved the RNA molecule to yield a  $^{32}$ P-labeled mono- and a  $^{32}$ P-labeled trinucleotide (not shown). Thus, the sequence of the tetranucleotide is 5'A-G-A-(U or C)3'. Based on these data, the complete structure of msDNA-Ec67 from *E. coli* Cl-1 is presented in Figure 7. Despite a lack of primary structural homology, msDNA-Ec67 displays all the unique features found in msDNAs from myxobacteria. These include a single-stranded DNA with a stem-and-loop structure, a single-stranded RNA with a stem-and-loop structure, a 2',5'-phosphodiester linkage between the RNA and DNA, and a DNA-RNA hybrid at their 3' ends. This hybrid structure was confirmed by demonstrating sensitivity of the RNA molecule to RNase H (not shown).

Cloning of the locus for msDNA-Ec67

In order to identify the DNA fragment which is responsible for msDNA synthesis in *E. coli* Cl-1, Southern blot hybridization was carried out with various restriction enzyme digests of total chromosomal DNA prepared from *E. coli* Cl-1, using msDNA-Ec67 labeled with AMV-RT (the same preparation as shown in lane 3, Figure 6) as a probe. The result is shown in Figure 9A. EcoRI (lane 1), HindIII (lane 2), BamHI (lane 3), PstI (lane 4) and BglII (lane 5) digestions showed single band hybridization signals corresponding to 11.6, 2.0, 2.2, 2.8 and 2.5 kilobase pairs (kb), respectively. The upper band appearing in the EcoRI digestion is due to incomplete digestion of the chromosomal DNA. Analysis of total chromosomal DNA prepared from *E. coli* Cl-1 by agarose gel electrophoresis revealed that the strain contains two plasmids of different size. However, neither plasmid hybridized with the <sup>32</sup>P- labeled probe, indicating that the fragments detected in Figure 9A are derived from chromosomal DNA. Furthermore, there is only one location for the msDNA-coding region on the chromosome, since various restriction enzyme digestions gave only one band of varying sizes. Similar results were observed for the msDNAs of myxobacteria (Yee *et al.*, 1984; Furuichi *et al.*, 1987; and Dhundale *et al.*, 1988).

The 11.6-kb EcoRI fragment and the 2.8-kb PstI fragment were each cloned into pUC9 (Yanisch-Perron *et al.*, 1985) and *E. coli* CL83 (a recA transductant of strain JM83), an msDNA-free K-12 strain (lane 1, Figure 9B), was transformed with the plasmids. Cells transformed with the 11.6-kb EcoRI clone (pCl-1E) were found to produce msDNA (lane 2, Figure 9B), whereas cells transformed with the 2.8-kb PstI clone (pCl-1P) failed to produce any detectable msDNA (lane 3, Figure 9B). A map of the 11.6-kb fragment is shown in Figure 10. Southern blot analysis of the fragment revealed that a 1.8-kb PstI - HindIII fragment hybridized with the msDNA probe. When the DNA sequence of this fragment was determined, a region identical to the sequence of the msDNA molecule was discovered. The DNA sequence corresponding to the sequence of msDNA is indicated by the enclosed box on the lower strand in Figure 11 and the orientation is from right to left. The location of this sequence is also indicated by a small arrow in Figure 10. As is the case for all other known myxobacterial msDNAs (Dhundale *et al.*, 1987; Furuichi *et al.*, 1987; and Dhundale *et al.*,

1988), a sequence identical to that of the RNA linked to msDNA (see Figure 7) was found downstream of the msDNA-coding region in opposite orientation and overlapping with that region by 7 bases. This sequence is indicated by the enclosed box on the upper strand in Figure 11 and the branched G residue is circled. Again, as in all the msDNAs found in myxobacteria, there is an inverted repeat comprised of a 13-base sequence immediately upstream of the branched G residue (residue 250 to 262; sequence a2 in Figure 11) and a sequence at the 3' end shown by an arrow in Figure 11 (residue 368 to 380; sequence a1). As a result of this inverted repeat, a putative longer primary RNA transcript beginning upstream of the RNA coding region and extending through the msDNA coding region would be able to self-anneal and form a stable secondary structure, which is proposed to serve as the primer as well as the template for msDNA synthesis (Dhundale *et al.*, 1987).

#### Existence of an essential gene for msDNA synthesis

The 2.8-kb PstI fragment (from PstI(a) to PstI(b) in Figure 10) was not able to synthesize msDNA. However, an overlapping 3.9-kb fragment from BalI (1.0 kb downstream of PstI(a); see Figure 10) to the following EcoRI site contains all the information required for synthesis of msDNA. This indicates that a region downstream of the PstI(b) site (Figure 10) is required for msDNA production. The nucleotide base sequence from this region revealed a long open reading frame (ORF) of 586 amino acid residues, starting with the initiation codon ATG at nucleotide 418 to 420 as shown in Figure 11. A distance of only 51 bases separates the initiation codon from the region which encodes msDNA. A putative Shine-Dalgarno sequence (GGA) can be found 10 bases upstream of the initiation codon. When the lacZ gene was fused in frame at the HindIII site (within the ORF) at amino acid residue-126,  $\beta$ -galactosidase activity was detected (not shown). Thus the region encompassing the ORF is indeed transcribed and the gene product encoded by the ORF is essential for msDNA synthesis. In a preliminary experiment, both msdRNA and the ORF appeared to be transcribed as the same transcription unit, since a deletion mutation removing the sequence from residue 1 to 181 blocked the expression of the lacZ gene fused at the HindIII site. A putative

promoter can be found in the deleted sequence as boxed in Figure 11. These -35 and -10 regions

probably serve as the promoter for both msdRNA synthesis and the ORF.

#### Sequence similarity with retroviral reverse transcriptases

When the amino acid sequence of the ORF was compared with known proteins, a striking similarity was found with retroviral RTs. In Figure 12, the ORF is compared with RTs from HIV (human immunodeficiency virus; Ratner *et al.*, 1985; and Johnson *et al.*, 1986), and HTLV1 (human T-cell leukemia virus type I; Seiki *et al.*, 1983; and Patarca *et al.*, 1984). The first domain (Asn-32 to Val-291) matches well with the RT domains of HIV and HTLV1. In particular, the sequences around the polymerase consensus "Asp-Asp" sequence (Toh *et al.*, 1983; and Geng *et al.*, 1985; boxed in Figures 11 and 12) are well conserved. Out of 260 amino acid residues in this domain, 44 and 38 residues are identical with HIV and HTLV1, respectively. Between HIV-RT and HTLV1-RT, there are 78 identical amino acid residues in this domain.

The *pol* gene of retroviruses is known to produce a protein consisting of RT and RNase H activities; the former at the amino-terminal and the latter at the carboxyl-terminal region of the *pol* gene product (Ratner *et al.*, 1985; Johnson *et al.*, 1986; Varmus, 1985; and Tanese *et al.*, 1988). These domains have been shown to be separated by a poorly conserved "tether" domain of approximately 160 to 190 amino acid residues (Ratner *et al.*, 1985; Johnson *et al.*, 1986). On the basis of the HIV sequence, the similarities (only identical amino acid residues) between HIV and HTLV1 are 29.5 and 16.8% for the RT domain and the tether domain, respectively. The similarities between HIV and msDNA are 16.9 and 10.3% for the RT domain and the tether domain, respectively. The similarities between HTLV1 and msDNA are 14.6 and 15.5% for the RT domain and the tether domain, respectively. These results indicate that in addition to the RT region, there are reasonable similarities in the tether domain between retroviruses and msDNA. An alignment of the RNase H domains also revealed that there are similarities between retroviruses and msDNA (15.7 and 17.4% with HIV and HTLV, respectively; see Figure 12). The similarity between HIV and HTLV1 in this region is 18.0%.

Cell extracts were prepared and assayed for the presence of RT activity associated with the production of msDNA as predicted from the amino acid homologies. Only the E. coli strain (C2110, polA) (Tanese et al., 1985; Tanese et al., 1986; E. coli strain C2110 (polA<sup>-</sup>) was a gift from M. Roth and S. Goff) harboring the plasmid, pCI-1EP5, containing the msDNA ORF displayed RT activity (Figure 13). The polA strain was used to eliminate high background activity in the RT assay due to DNA polymerase I. No RT activity was detected in extracts containing the vector plasmid alone, or when the template-primer (poly rC-dG) was absent from the reaction mix (Figure 13). It is interesting to note that the PstI(b) site is located at amino acid residue-430, which is between the tether domain and the RNase H domain. A plasmid lacking sequences downstream of the PstI(b) site did not produce msDNA. This suggests that the RNase H domain may be essential for msDNA synthesis, or alternatively that PstI disruption may result in inactivation of RT.

In addition to the similarity between msDNA-Ec67 RT and retroviral RT, there is an interesting similarity between msDNA and retroviruses; DNA synthesis starts at a site upstream of the RT-RNase H gene, and the orientation of DNA synthesis is opposite to the direction of transcription of the RT-RNase H gene. In the case of retroviruses, tRNAs are recruited from the cell for the priming reaction (Weiss et al., 1985), whereas for msDNA an RNA transcript serving as, template also serves as a primer by self-annealing to form a stable secondary structure (Dhundale et al., 1987; Furuichi et al., 1987).

#### Origin of the E. coli Reverse Transcriptase

At present the relationship between msDNA and retroviruses is an open question. It is possible that the study of msDNA may shed light on the question of the origin and evolution of retroviruses. It is an intriguing question to consider why some of the clinical E. coli strains, isolated from human patients produce msDNA. Our preliminary data indicate that msDNAs produced by four independent E. coli strains, isolated from urinary track infections, share little homology. This suggests that there may be enormously large numbers of species of msDNA in E. coli. In contrast to

msDNAs found in E. coli, msDNA-Mx162 from M. xanthus is highly conserved, since nine

independent M. xanthus strains isolated from various sites have msDNA which hybridizes with the original msDNA-Mx162 (Dhundale et al., 1985). Furthermore, msDNA from another myxobacterium, S. aurantiaca (msDNA-Sa163; Furuichi et al., 1987), also shows a high degree of homology to msDNA-Mx162 (Furuichi et al., 1987).

5

Several lines of evidence suggest that the RT gene found in the E. coli strain Cl-1 is not likely to have originated in E. coli, but rather was recently acquired from some other source. For example, only about 4% of E. coli strains tested were found to produce msDNA. In addition, the RT gene from strain Cl-1 does not cross hybridize to chromosomal DNA from four other E. coli strains which produce msDNA molecules, indicating that there is extensive diversity among these RT genes.

10

In contrast, a DNA fragment from the E. coli-K-12 sigma factor gene can hybridize to chromosomal DNA from all five msDNA producing, E. coli strains, indicating the conserved nature of sigma factors. An analysis of the E. coli RT gene indicates that the codon usage for this gene is remarkably different from most E. coli proteins. In particular, AGA and AGG, the least frequently (2.7%) used codons for arginine among 199 E. coli genes (Maruyama et al., 1986), occurs at a frequency of 64.5% in the E. coli RT gene. Similarly, CUG is the most commonly used codon for leucine (61.3%; Maruyama et al., 1986) in E. coli genes, while its prevalence in the RT gene is only 9.1%. The AT base pair content of the E. coli RT gene was calculated to be 67.6%, which is substantially higher than the AT content of the E. coli genome (45%; Fasman, 1976). The AT contents of HIV and HTLV1 RT genes are 62.1% and 47.8%, respectively. These facts pose an intriguing question as to how and when the RT gene, as well as the msDNA coding region, were integrated into the genome of the clinical strain.

20

There are many questions to be answered, including (a) are there any particles associated with msDNA, (b) is the msDNA region transposable like the Ty element of yeast (Boeke et al., 1985; Eichinger et al., 1988), (c) can the element responsible for the production of msDNA be transferred from cell to cell, (d) can a RT from one strain (E. coli or myxobacteria) complement the production of msDNA of other strains, (e) does the promoter for the RNA transcript have any

25

similarities to the retroviral LTR, (f) are there any specific integration sites for the msDNA element

on the *E. coli* chromosome, (g) why is the branched G residue conserved, (h) is there an enzyme responsible for priming DNA synthesis at the 2'-OH position of the rG residue, (i) why and how does msDNA synthesis stop at one distinct site on the RNA template, and (j) how different biochemically are the msDNA RTs from retroviral RTs?

5           The existence of reverse transcriptase in prokaryotes, previously speculated upon (Dhundale et al., 1987), is now evident. This fact raises intriguing questions concerning possible roles of this enzyme in the prokaryotes other than a role in msDNA production. Recently we also found that *M. xanthus*, in which msDNA was originally discovered, has a long ORF in the same manner as found for msDNA-Ec67. This ORF has a high degree of similarity to the *E. coli* RT. Since eight  
10 independent isolates of *M. xanthus* produce homologous msDNA, the *M. xanthus* RT is likely to have been acquired at a very early stage of its evolution in contrast to the *E. coli* RT. The determination of the structures of both *M. xanthus* and other *E. coli* RTs will shed light on the key question of the origin of RT and its role in prokaryotes.

          An important embodiment of the invention relates to the discovery of msDNA-producing retron elements in a number of diverse bacterial groups. Thus, retron elements appear to be widely prevalent, at least amongst the purple bacteria or proteobacteria including *Proteus*, *Klebsiella* and *Salmonella* of the gamma subdivision; *Rhizobium* and *Bradyrhizobium* from the alpha subdivision; and *Nannocystis* (a myxobacterium) from the delta subdivisions. These are representatives of the three of the four major subdivisions of the purple bacteria of proteobacteria.  
20 As shown above the retron-encoded RT is responsible for the synthesis of msDNAs.

          The retron elements were discovered by detecting the presence of msDNA by one of two classic methods: the so-called "RT extension method", described by Lampson, B.C., M. Inouye and S. Inouye, 1991. Survey of multicopy single-stranded DNAs and reverse transcriptase genes among natural isolates of *Myxococcus xanthus*. J. Bacteriol. 173:5363-5370 and in Lampson, B.C.,  
25 M. Viswanathan, M. Inouye and S. Inouye, 1990. Reverse transcriptase from *Escherichia coli* exists as a complex with msDNA and is able to synthesize double-stranded DNA. J. Biol. Chem. 265:8490-

8496 or polyacrylamide gel electrophoresis of a chromosomal DNA extract followed by staining with



ethidium bromide as described by Yee, T., T. Furuichi, S. Inouye, 1984. Multicopy Single-Stranded DNA Isolated from a Gram-Negative Bacterium, Myxococcus xanthus. Cell, Vol. 38, 203-209. Both of these publications are incorporated herein by reference. Both methods provide a reliable, convenient and conventional protocol for screening of bacteria for the presence of retron-encoded RT and msDNAs.

In accordance with the RT extension method, the DNA portion of msDNA is specifically  $^{32}\text{P}$  radio labeled. Radio labeled from a total RNA preparation extracted from each bacteria strain to be screened. Twenty or more isolates of proteus mirabilia, Klebsiella pneumoniae, Salmonella species, rhizobial species, and enterococcal species were screened by this method. Low-molecular-weight bands (Fig. 20) indicated the presence of small labeled DNAs after polyacrylamide gel electrophoresis and autoradiography of the labeling reaction mixes. In addition, half of each labeling reaction mix was also treated with RNase A, causing a shift to a faster-migrating band, indicating that the labeled DNA is also associated with RNA. This is hallmark of the msDNA molecule as discussed above. Four of the 23 P. mirabilia isolates screened produced msDNA, while only 1 of 21 K. pneumoniae isolates and 4 of 70 Salmonella isolates screened produced msDNA. msDNA was detected in any of the 30 or so enterococcal strains screened by this method. It was concluded that the bacterial genera which contain msDNA producing retron elements are representatives of three of the four major subdivision of the purple bacteria or Proteobacteria, as described above.

In accordance with this embodiment of the invention, it is noteworthy that the discovery of msDNA extends for the first time the distribution of retron-elements to a new phylogenetic division of the purple bacteria, namely, the alpha subdivision. A collection of 63 rhizobial isolates (shown in Table 1) were screened for the presence of msDNA by the RT extension method. Among the 63 isolates, msDNA were detected in 10 (16% - Fig. 20 and Fig. 21). However, all 10 positive isolates give strong, clearly labeled bands with a typical shaft of a fast-migrating band after treatment with RNase A, indicating the presence of RNA and DNA in the labeled molecule.

The 10 retron-encoding rhizobial strains include both fast growing (rhizobium) and slow-growing

(Bradyrhizobium) rhizobia.

The RT extension method comprises treating a preparation of total RNA, extracted from a bacterial strain to be tested, with RT from a suitable source in the presence of the deoxynucleotides dATP, dTTP, dGTP and dCTP, one of which is radiolabeled, e.g., [ $\alpha$ - $^{32}$ P] dCTP, electrophoresing the treated RNA preparation on a polyacrylamide gel and determining initially the presence or absence of msDNA in the bacterium of interest by detecting a band of radiolabeled DNA corresponding to the single-stranded DNA of msDNA. Typical examples of suitable sources of RT are avian myeloblastosis virus (AMV) and Moloney murine leukemia virus (Mo-MLV). Conceivably, the test could be automated.

Total RNA samples, which contain msDNA if present in the bacterium, are extracted from the bacterial strain of interest and prepared for RT extension as follows. Total RNA, prepared from a 5-ml culture from the bacterial strain, is added to 50  $\mu$ l of a reaction mixture containing: 50 mM tris-HCl (pH 8.3); 6 mM  $MgCl_2$ ; 40 mM KCl; 5 mM DTT; 1  $\mu$ M dATP, dTTP and dGTP; 0.04  $\mu$ M dCTP; 0.2  $\mu$ M [ $\alpha$ - $^{32}$ P] dCTP; and 10 units of AMV-RT (Boehringer Mannheim). The reaction mixture is incubated at 37°C for 30 minutes, then extracted with 50  $\mu$ l of phenolchloroform (1:1) and precipitated with ethanol. The samples are subjected to electrophoresis on a 4% acrylamide - 8 M urea gel with appropriate nucleotide size markers, e.g., the Klenow fragment of DNA polymerase I. If the labeled sample is digested with ribonuclease (RNase) A before it is placed on the gel, a single band corresponding to single-stranded DNA is detected, which is indicative of the presence of msDNA. An aliquot from each labeling reaction mixture is treated with 5  $\mu$ g of RNase for 10 minutes at 37°C just prior to electrophoresis to detect in the gel a shift to a faster - migrating species, indicating that each labeled DNA is also associated with RNA, which is the hallmark of the msDNA molecule.

Low-molecular weight bands in the gel indicate the presence of small labeled DNAs after polyacrylamide gel electrophoresis and autoradiography of the labeling reaction mixtures.

Multiple bands observed in some of the lanes of the gel even after RNase treatment may be due to incomplete extension by RT during the labeling reaction, or, alternatively, multiple forms or species of msDNA may exist in a given bacterium.

The Yee method for screening bacteria for the presence of retrons which synthesize msDNAs involves purifying by a conventional phenol extraction procedure total chromosomal DNA from the desired bacteria to be screened, electrophoresis on a five percent preparation acrylamide gel and checking for a satellite band. The major satellite band is cut out to extract the material in the band to quantitate the material in the satellite band. Total chromosomal DNA is subjected to acrylamide gel electrophoresis, the gel is stained with a ethidium bromide and densitometric scanning is employed to quantitate the satellite DNA against the pBR322 standard. The method is described in better details in Yee cited above.

A collection of rhizobial isolates from the United States Department of Agriculture (USDA) Beltsville Rhizobium Culture Collection are screened for the presence of msDNA by the RT extension method. This collection represents isolates at different times, from different legume hosts and from different geographic locations. msDNAs are detected in 10 isolates. All 10 positive isolates give strong, clearly labeled bands of DNA, with a typical shift to a fast-migrating band after treatment with RNase A, indicating the presence of RNA and DNA in the labeled molecule. The 10 retron-encoding rhizobial strains include both fast-growing (Rhizobium) and slow-growing (Bradyrhizobium) rhizobia as follows: Rhizobium sp. (Acacia) 3002 and 3838, Bradyrhizobium sp. (Aeschynomene) 3516, Bradyrhizobium sp. (Albizia) 3004, Bradyrhizobium sp. (Erythrina) 3242, Rhizobium loti 3468 and 3503, Rhizobium trifolii 2048 and 2065 and Bradyrhizobium sp. (Vigna) 3447. See Figure 21

Total DNA from each of eight msDNA-producing strains clearly cross-hybridizes with a nod YAB (1.6 - kb Eco RI fragment) gene probe derived from Bradyrhizobium japonicum, confirming that these strains are members of the Rhizobiaceae.

In view of the diversity of retron elements in prokaryotic populations, it is not excluded that msDNA synthesizing retrons would be found in bacteria living in alkaline environments, such as in alkaline environments: Plectonema nostocorum, Flavobacterium spp.

Agrobacterium spp. Bacillus spp. Ectothiorhodospira spp.; in acidic environments: Thiobacillus thermophilica and thiooxidans, Thermoplasma acidophilus, Sulfolobus acidocaldarius, Cuanidium

caldarius, Bacillus acidocaldarius; in very high temperature environment (thermophilic): Sulfolobus  
acaidocaldarius, Caldariella acidophila, Thermus aquaticus; in very low temperature (psychrotrophic):  
Vibrio marinus, Pseudomonas spp., Cytophaga spp., Flavobacterium spp.; in high salt environments  
(halophilic): Halobacterium cutirubrum and salinarium, Halococcus morrhuae, Danaliella viridis; in  
5 high barometric pressure (like deep sea - barophilic), which are believed to inhibit the gut of ocean  
bottom dwelling fish. By using one of the two screening tests identified above, one skilled in the art  
will readily determine whether any one of these bacteria contain retrons synthesizing msDNA. This  
may be particularly interesting for making evolutionary comparisons between homologous RT genes  
present in distantly related phytogenic strains.

10 A representative number of amino acid sequences of representative RTs were analyzed  
to determine similarities and differences. The following observations were made. The amino acid  
sequences of these bacterial RTs are shown in Figure 14. The individual nucleotide and amino acid  
sequences for each of the RTs are shown in Figures 2, 11 and 15 through 19.

15 From a comparison of these sequences, it is noted that there are 61 conserved positions  
in the RT domains as indicated by solid dots at the bottom of the sequences in Figure 14. It is further  
noted that all bacterial RTs possess the YXDD sequence. Several other residues are conserved  
including the LPQS sequence that is especially common in retroviral reverse transcriptases. The RT  
domains are divided into seven subdomains. For each subdomain, the consensus sequences for the  
seven bacterial RTs can be established, as shown at the bottom of the sequences in Figure 14. There  
20 are 18 extra residues (except 26 residues for RT-Ec67) between subdomains 2 and 3, in which there  
is a reasonably good consensus sequence.

It has been noted that the RTs of the present invention possess a number of common  
conserved sequences of nucleotides and amino acid residues.

25 The most common conserved sequence of amino acid residues noted is as follows:  
tyrosine, alanine or cysteine and two aspartic acid residues. This conserved sequence, common to all  
RTs of the present invention, is also known as the YXDD sequence.  
*as shown in Seq. ID No. 4350*

A second conserved sequence of amino acid residues noted is as follows: serine, x which is a hydrophobic residue selected from the group consisting of valine, phenylalanine leucine and isoleucine,  $x_1$  which is a polar residue selected from the group consisting of threonine, asparagine, lysine and serine and  $x_2$  which is a hydrophobic residue selected from the group consisting of tryptophan, phenylalanine and alanine. *as shown in Seq. ID No. 4451*

A third conserved sequence of amino acid residues noted is as follows: asparagine, x which is a hydrophobic residue selected from the group consisting of alanine, leucine and phenylalanine and  $x_1$  which is a hydrophobic residue selected from the group consisting of leucine, valine and isoleucine. *as shown in Seq. ID No. 4452*

A fourth conserved sequence of amino acid residues further noted is as follows: x which is a polar residue selected from the group consisting of arginine, glutamic acid, lysine, valine and glutamine, a second residue which is valine, a third residue which is threonine and a fourth residue which is glycine. *as shown in Seq. ID No. 4453*

These conserved sequences are only a portion of the total number of common sequences of the RTs. For other conserved sequences held in common by the bacterial RTs reference is made to Figure 14.

The RTs of the other groups of bacteria described herein as capable of synthesizing msDNAs are likewise believed to have a similar profile of conserved nucleic acid and amino acid residue sequence similarities as shown in Figure 14 and discussed above. This observation also applies to the genus Nannocystis.

In accordance with the invention, it is contemplated that prokaryotic reverse transcriptase, which is essential for msDNA synthesis, may be responsible for host cell parasitic or selfish DNA synthesis. Additionally, it is thought that the prokaryotic reverse transcriptase molecule may be essential for synthesis of biological messengers and nucleic acid enzymes.

The msDNAs synthesized by the reverse transcriptase disclosed herein possess a highly stable RNA; it is capable of self-annealing and may serve as the primer and template for msDNA synthesis. The reverse transcriptases (RTs) disclosed herein may be used as diagnostic agents. It is

also contemplated that the RTs of the invention can synthesize msDNAs which will contain specific selected DNA fragments that can hybridize with complementary ssDNA, or otherwise identify ssDNAs, sought for, thus being useful as probes.

The possibility for the msDNAs to behave like restriction enzymes (or have restriction-like enzyme activity) in being capable of cleaving DNAs, or cut off a segment of itself, cannot be excluded.

The following examples are provided for purposes of illustration only and are not to be viewed as a limitation of the scope of the invention. The following examples are illustrative of bacterial isolates screened and identified to contain msDNA by way of the present invention.

#### EXAMPLE 1

One of the rhizobial strains, Rhizobium trifolii USDA 2065 is identified as containing msDNA by the RT extension method by which msDNA from total RNA is specifically labeled with  $^{32}\text{P}$  as follows.

Total RNA from a 5-ml culture of R. trifolii 2065 is added to a 50  $\mu\text{l}$  reaction mixture containing: 50 mM tris-HCl (pH 8.3); 6 mM  $\text{Mg Cl}_2$ ; 40 mM KCl; 5 mM DTT; 1  $\mu\text{M}$  dATP, dTTP and dGTP; 0.04  $\mu\text{M}$  CTP; 0.2  $\mu\text{M}$  [ $\alpha\text{-}^{32}\text{P}$ ] dCTP; and 10 units of AMV-RT (Boehringer Mannheim). The reaction mixture is incubated at 37°C for 30 minutes, then extracted with 50  $\mu\text{l}$  of phenolchloroform (1:1) and precipitated with ethanol. The samples are subjected to electrophoresis on a 4% acrylamide-8 M urea gel with appropriate nucleotide size markers, such as the Msp I digest of pBR322 end-labeled with [ $\alpha\text{-}^{32}\text{P}$ ] dCTP and the Klenow fragment of DNA polymerase I. An aliquot of the reaction mixture containing R. trifolii RNA is treated with 5  $\mu\text{g}$  of RNase for 10 minutes at 37°C prior to electrophoresis to detect in the gel a shift to a faster-migrating species, which indicates that the  $^{32}\text{P}$ -labeled DNA extended by RT is also associated with RNA, which clearly demonstrates the presence of msDNA.

Low-molecular weight bands in the gel indicate the presence of small  $^{32}\text{P}$ -labeled DNA after polyacrylamide gel electrophoresis and autoradiography. The labeled DNA is indicative of the presence of msDNA.

### EXAMPLE 2

5 By the method described above in Example 1, (a) Proteus mirabilis 1174b is found to synthesize msDNA by the retrons containing the RT; (b) Klebsiella pneumoniae 912b is found to synthesize msDNA by RT; (c) Salmonella sp. strain SARB-3 is found to synthesize msDNA by the retrons containing the by the retrons containing the RT; (d) Nannocystis exedens Nael is found to synthesize msDNA by RT; (e) Bradyrhizobium spp. 3447, 3516 and 3004 are also found to synthesize msDNA by the retrons containing the RT.

The following method, exemplified for E. coli, for the isolation and purification of bacterial RT is applicable to bacteria which are screened as positive for the presence of msDNA by the RT extension in vitro method.

### EXAMPLE 3

#### Isolation and Purification of Bacterial Reverse Transcriptase.

The following is a description of a convenient method for isolating and purifying a bacterial RT.

From 10 liters of a stationary phase culture of E. coli strain C2110 harboring plasmid pCl-1EP5b, cells are harvested, washed in 50 mM Tris (pH 8.0), and resuspended in lysozyme buffer (50 mM Tris (pH 7.5), 10% sucrose, 0.3 M NaCl, 1 mM EDTA, 1 mM phenylmethylsulfonyl fluoride). Fresh lysozyme is added to a final concentration of 2 mg/ml. The suspension is incubated on ice for 15 minutes followed by a quick freeze at  $-70^{\circ}\text{C}$ , then thawed on ice. Lysis is enhanced by the addition of 2 volumes of buffer M (50 mM Tris (pH 7.0), 1 mM dithiothreitol, 0.2% Nonidet P-40,

10% glycerol, and 25 mM NaCl) followed by incubation on ice, then a quick freeze-thaw. A cleared lysate is obtained by centrifugation at 38,000 rpm in a 50Ti rotor for 30 minutes. The cleared lysate is fractionated by ammonium sulfate precipitation (0-50%, 50-70% and 70-90%), followed by dialysis overnight (4°C) for each fraction against buffer M. Ammonium sulfate fractions, 50-70% and 70-90%, show RT activity and are pooled, then applied to a DEAE-column (2.5 x 50 cm; DE52 Whatman) equilibrated with buffer M. The DE52 column is washed, and RT activity is eluted from the column at a range of 300 to 350 mM NaCl. The DE52 fractions showing RT activity are pooled, concentrated by membrane ultrafiltration (Amicon) and then loaded onto a Sephacryl S-300 column (Pharmacia LKB Biotechnology Inc., 1.5 x 75 cm) equilibrated with buffer M. The column is developed with the same buffer. Again, fractions from the S-300 column having RT activity are pooled and concentrated, and 0.7 ml is loaded onto a 16-30% glycerol density gradient. The glycerol gradients are set up and run as described previously (Viswanathan *et al.*, 1989). The purified Ec67.RT (fractions 7, 8 and 9) is stored as separate glycerol fractions at -20°C.

When this protocol is applied to the msDNA bacterial synthesizing strains, the respective RTs are isolated and identified as shown above.

Another convenient method for isolating and purifying reverse transcriptase is published in Lampson B.C., S. Inouye and M. Inouye, "msDNA of Bacteria", Progress in Nucleic Acid Research and Molecular Biology, Vol. 40, pages 1 *et seq.*

The invention has been described in detail with particular reference to the above embodiments. It will be understood, however, that variations and modifications can be affected within the spirit and scope of the invention.



CLAIMS

We claim:

1. An isolated and purified bacterial reverse transcriptase (RT) which is capable of synthesizing msDNA, which RT comprises a conserved sequence of amino acid residues as follows: tyrosine, x which is alanine or cysteine, and two aspartic acid residues.

2. The bacterial RT of claim 1 which comprises a second conserved sequence of amino acid residues as follows: serine, x which is a hydrophobic residue selected from the group consisting of valine, phenylalanine, leucine and isoleucine,  $x_1$  which is a polar residue selected from the group consisting of threonine, asparagine, lysine and serine and  $x_2$  which is a hydrophobic residue selected from the group consisting of tryptophan, phenylalanine and alanine.

3. The bacterial RT of claim 2 which comprises a third conserved sequence of amino acid residues as follows: asparagine, x which is a hydrophobic residue selected from the group consisting of alanine, leucine and phenylalanine and  $x_1$  which is a hydrophobic residue selected from the group consisting of leucine, valine and isoleucine.

4. The bacterial RT of claim 3 which comprises a fourth conserved sequence of amino acid residues as follows: x which is a polar residue selected from the group consisting of arginine, glutamic acid, lysine, valine and glutamine, a second residue which is valine, a third residue which is threonine and a fourth residue which is glycine.

5. The bacterial RT of claim 1 which has the common subdomains 1 through 7 shown in Table 5.

6. The bacterial RT of claim 1 wherein the conserved sequence is located in subdomain 5 shown in Table 5.

7. The bacterial RT of claim 6 which has a total of 61 conserved amino acid residues.

8. An isolated and purified bacterial RT which comprises a sequence of amino acid residues shown in Figure 14.

9. An isolated and purified bacterial RT from a bacterium which is capable of synthesizing an msDNA as determined by the reverse transcriptase extension in vitro screening test, which indicates the presence or absence of msDNA in the bacterium.

10. The bacterial RT of claim 9 wherein the bacterium is selected from the group of genera consisting of Myxococcus, Escherichia, Proteus, Klebsiella, Flexabacter, Cytophaga, Stigmatella, Salmonella, Nannocystis, Rhizobium and Bradyrhizobium.

11. The bacterial RT of claim 10 wherein the in vitro screening test for determining the presence or absence of msDNA in the bacterium comprises treating a preparation of total RNA extracted from the bacterium with a reverse transcriptase (RT) in the presence of a radiolabeled deoxynucleotide, which RT, when msDNA is present in the total RNA of the bacterium, utilizes the DNA portion of the msDNA as a primer and the RNA portion of the msDNA as a template for radiolabeling the DNA portion of the msDNA, electrophoresing the treated RNA preparation and determining the presence of msDNA in the bacterium by detecting a band of radiolabeled DNA, said band being indicative of the presence of msDNA in the bacterium.

## REFERENCES

- LAW OFFICES  
WEISER & ASSOCIATES**  
SUITE 500  
230 SO. FIFTEENTH ST.  
PHILADELPHIA, PA 19102  
(215) 875-8383  
TELECOPIER (215) 875-8394

Vieira J., Messing J., Gene, 19, 259-268 (1982).

Visawanathan M., Inouye M., Inouye S., J. Biol. Chem., 264, 13665-13671 (1989).

Voytas D.F., Ausbel F.M., Nature, 336, 242-244 (1988).

Weiss N., Teich H., Varmus H., Coffin J., RNA Tumor Viruses, Vol. 2, Cold Spring Harbor Laboratory (1985).

Yee T., Furuichi T., Inouye S., Inouye M., Cell, 38, 203-209 (1984).

Yuki S., Ishimaru S., Inouye S., Saigo K., Nucl. Acid Res., 14, 3017-3020 (1986).

03606431.030397



255060-1000000

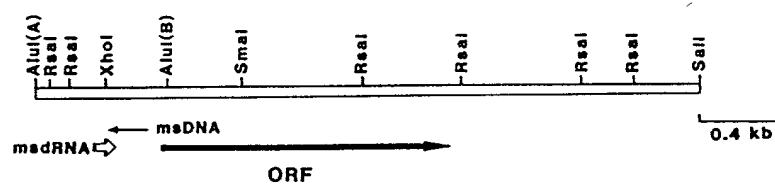


FIGURE 1

[illegible]

FIGURE 2



HIV	RT	VKLPGMDGPKVKQ	WPLTEEKIKALVEICTEMEKEGKISKIGPENPYNTPVFAIKKDKSTKWR	
HTLV1	RT	RPWARTPPKAPRNQ	PVPFKPERLQALQHLVRKALEACHIEPYTG	PCNNPVFPVKKA NGTWR
Ec-67	RT	NVLYRGDSNQYQTFTI	PKPKGKVRTISAPTDLR	KDIQRICDLLSDCRDEIFAIRKI SNNSYS
Mx-162	RT	AFHREVDTIATHYVSWIT	PKRDGSKRTITSPKPEL KAAQR	WVLS NVV ERLP VHGA
			o o o o	o o o o x o x x o o x o
HIV	RT	KLVDFRELNKRTQDFWEVQLGIPHAGLKKK	KSIVTLVDVGDAYFVSPVLDEDFRKYT	A
HTLV1	RT	FHDLRATNSLTIDLSSSPGPDLSLPTLAHLQTLIDLRDAFFQILPKQCFPYF		A
Ec-67	RT	FGFE RCKSIILNAYKHRCQIILNIDLKDFEFSEFNFRGVRG	YFLS NQDF	L
Mx-162	RT	HGCV AGRSILNALAHQCADVVVVKVDLKDFFPSVTVRRVKGLLRKGLREGTSTILLSLSTEAP		
			o o o o x x o o	o o o o x x x o o x x x
HIV	RT	FTIP SINNETPGIRYQYNVLPQGWKGSIPAIFGS	SMTKILEPKFKQNPDIIVYQ	YMDLLVYG
HTLV1	RT	FTVP QQCNYPGCTRYAWKVLQGVFKNSPTLFEM	QLAHILQPIRQAFQCTILQYMDLLILA	
Ec-67	RT	LN FVVATTLAKAACYN GTLPQGSPCSPISILNGLICINHMDRLAKLAKKY	GCTYSRYADITI	
Mx-162	RT	REAVQFRGKLHLVAKGP RALPGAGTSPGITNALCLKLKRLSALAKRL	GFTYTRYADLLTF	
			••••• o •• x o	x x o o x o x o x o o o o x
HIV	RT	S DLEIQHRTKIEELRQHLLRWGLTTP	DKKHQKEP PFLWMGYELHPDKWTVQPIVLPE	KD
HTLV1	RT	S PSHEDLLLSEATHASLISHGLPV	ENKQTQPTGTIKFLGQIISPNHLTYDAVPTVPI	RS
Ec-67	RT	STINKNTPFLEMATVQPEGVVLGKVLKEIENSGFEINSDKSLRTLYKTSRQEV	GLTVNRIVNID	
Mx-162	RT	SWTKAKQPKPRRTQRPVAVLLSRVQEVVAEGFRVHPDKTRVARKGTGRQVT	GLVVNAAGKDA	
			o o o o x o	x o o o x o x o x x
HIV	RT	SWTVNDIQKLVGKLINWASQIYP		
HTLV1	RT	RWALPELQALLGEIQWVSKGTP		
Ec-67	RT	RCYYKKTRALAHALYRTCE YK		
Mx-162	RT	PAARVPDVRQDLRAAIHN RK		
			x	

FIGURE 3

# A

Mx-162 18 PTELTAPSSDAAAKREARRLAHEALLVRKAIDEAGGADDWVQAQLVSKGLAVEDLD-FSSASEKDKA-WKEKK 91  
 Mo-MLV 1070 PDPDMTRVTNSPSLQAHLQALYLVQHEVW-RPL-AAAYQEQ-LDRPVVPHPYRVGDTVWVRRHQTKNLEPRWKGPPY 1142

Mx 162 92 KAEATERRALKRQAHEAW-KATHVGHGAGVHWAEDRL 128  
 Mo-MLV 1143 TVLLTTPTALKVDGIAAWIHAHVKAADPGGG-PSSRL 1179

# B

Mx-162 411 GKDAPAAARVPRDVVRQLRAAIHNRKKGKPGREGESLEQLKGMAAFIHMTD-PAKGRAFLAQLTELESTASAAAPQAE 485  
 HIV 396 GKEGHSARQCR-APR--RQGC--WKCGKPGHINTNCPD-R-QAGFLGLGPGWKKPRNFPVAQVPQ-GLTPTAPP 461

Figure 4. Sequence Similarity of the msDNA-Mx162 Reverse Transcriptase with Other Retroelements

(A) Sequence similarity of the region from residues 18 to 128 of the msDNA Mx162 RT (see Figure 2) with a carboxy-terminal region of integration protein of Moloney murine leukemia virus (M-MuLV) (residues 1070 to 1179, Shinnick et al., 1981)  
 (B) Comparison of the sequence from residues 411 to 485 of the msDNA-Mx162 RT (see Figure 2) with the sequence from residues 396 to 461 of the gag protein of human immunodeficiency virus (HIV, Ratner et al., 1985)

FIGURE 4

## A

Mx-162	304	GP-RALPQGAPTSPGITNALCLKDKRLSALAKRL-GFTYTRYADLTLF-SWTAKQPKPRRTQRPVAVL	371
Ec-67	159	YN-GTLPQGSPSPSIINLICNMHMLAKLAKKY-GCTYSRYADDITI-STNKWTFPLEMATYQPEGVVL	226
Ec-86	130	YK-NLLPQGAPSSPKLANLICSXLDYRIQGYAGSR-GLIYTRYADLTL-SAQSMKKVVKARDFLFSIIPS	197
HIV	311	YQYNVLPQGKGSSPAIFQS---SMTRILEPFFKKQNPDIVIYQYMDLLVYGS-DLETGQHRTKIEELRQHLL	377
HTLV1	150	YAWKVLPGGFKNSPTLFEM---QLAHILQPIRQAFQCTILQYMDILLAS--PSHEDLLLSEATMASLI	215
Mo-MLV	303	LWTWRLPGGFKNSPTLFDE---ALHRDLADFRIQHPDLILQYVDDLALLAA-TSELDCCQG-TRALL-QTL	367
RSV	141	FQWKVLPGMTCSP TICQL---VVGQVLEPLRLKHPSLCMLHYMDILLAA---SSHDLGEEAAGEEVI-STL	205
BLV	122	FAWRVLPGGFINSPALFER---ALQEPLRQVSAAFSQSLLVSYMDILLIYAS--PTEEQRSQCYQALA-ARL	186
Mt. plasmid	288	IATNGVPGGASTSCGLATYNVL-----ELFLRY--DELIMADGIL-CRQDPSTPDFSVEEAGVQEP	348
17.6	339	YEYLRMPFGLKNAP-ATFQRCHN-DI---LRPLLNKHC-LVYLDIIIVFS-TSLDEHLQSLGLVFE--KL	399
GYPSY	284	YEFCLRPGLRNASSIFQR---ALDDV---LREQI-GKICYVVDVVIIFS--ENESDHVRHIDTVLK-CL	344
Copia	1032	CKLNKAIYGLKQAARCFR-CIYI---LDKGNINENIYV-LLVVDVVIAT--GDMTRMNFKRYLME-KF	1112
Tal-3	990	CLLKKSLYGLKQSPROWNA-CVYV-KQVSE-QEHLYL---LLVVDMLIAG--KSKSEINKVKEQLSM-EF	1069
Ty912	948	IRLKKSLYELKQS-CANWYE--EVRG-WSCVFKNQV-TICLFVDHMLVFS--KNLNSNKRIIEKLKM-QY	1023

# B

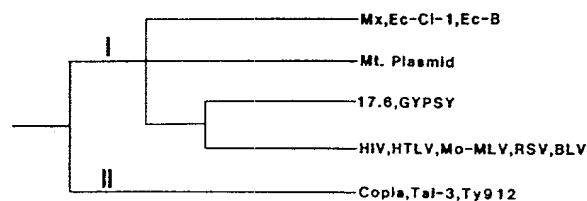
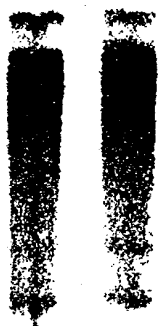


FIGURE 5

S 1 2 3 4



030397

FIGURE 6

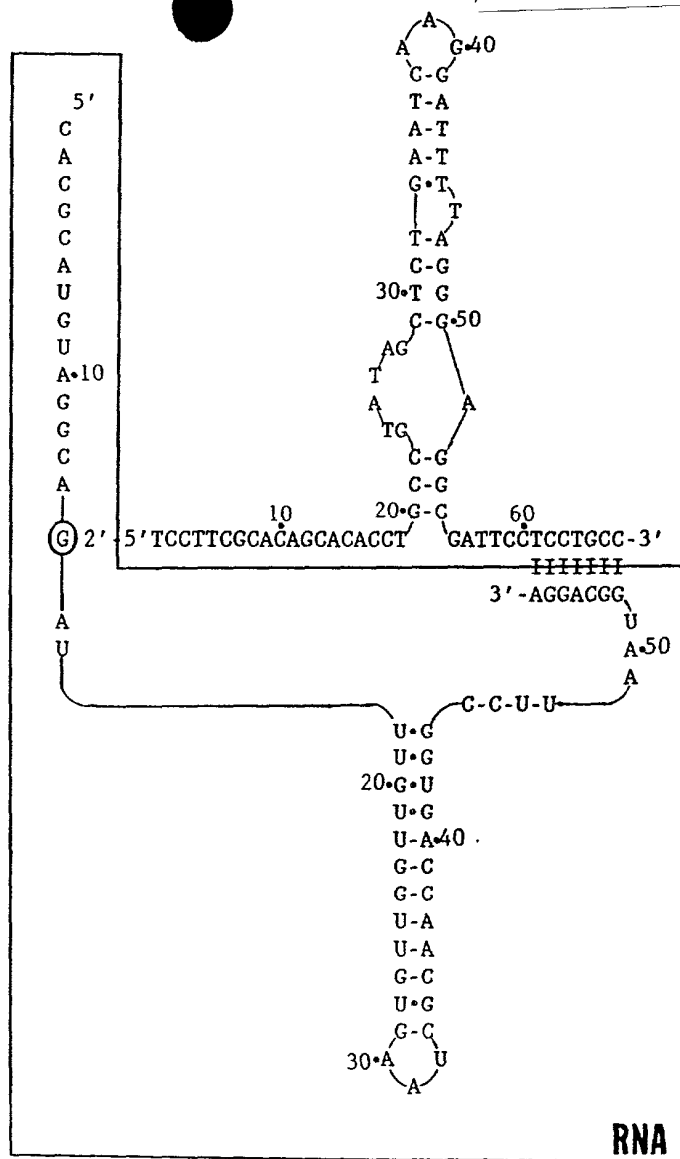


FIGURE 7

0806031 030397

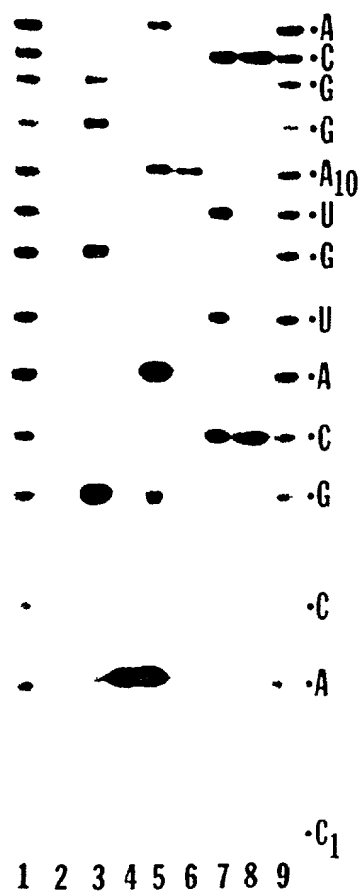
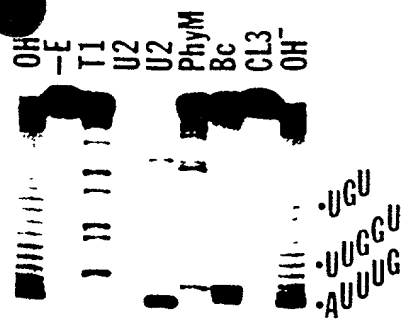


FIGURE 8

**A**

1 2 3 4 5

23.0-

9.4-

6.6-

4.4-

2.3-

2.0-

268950-1-030397

**B**

S 1 2 3

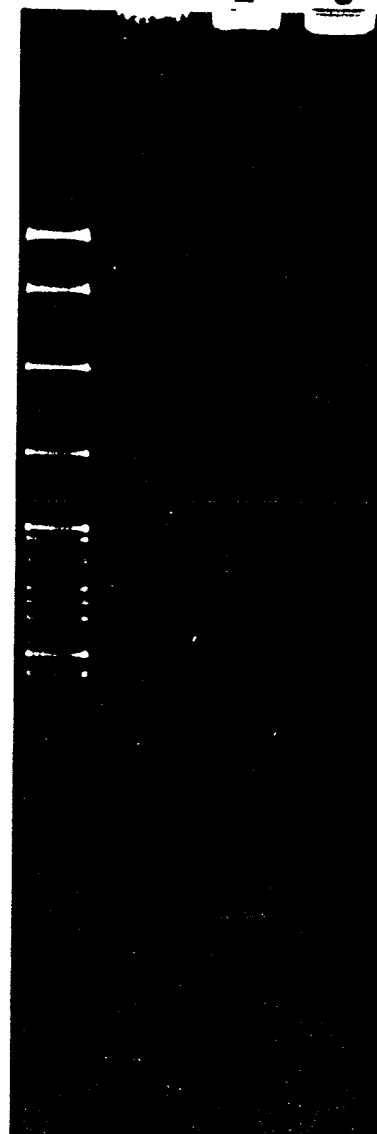


FIGURE 9

FIGU

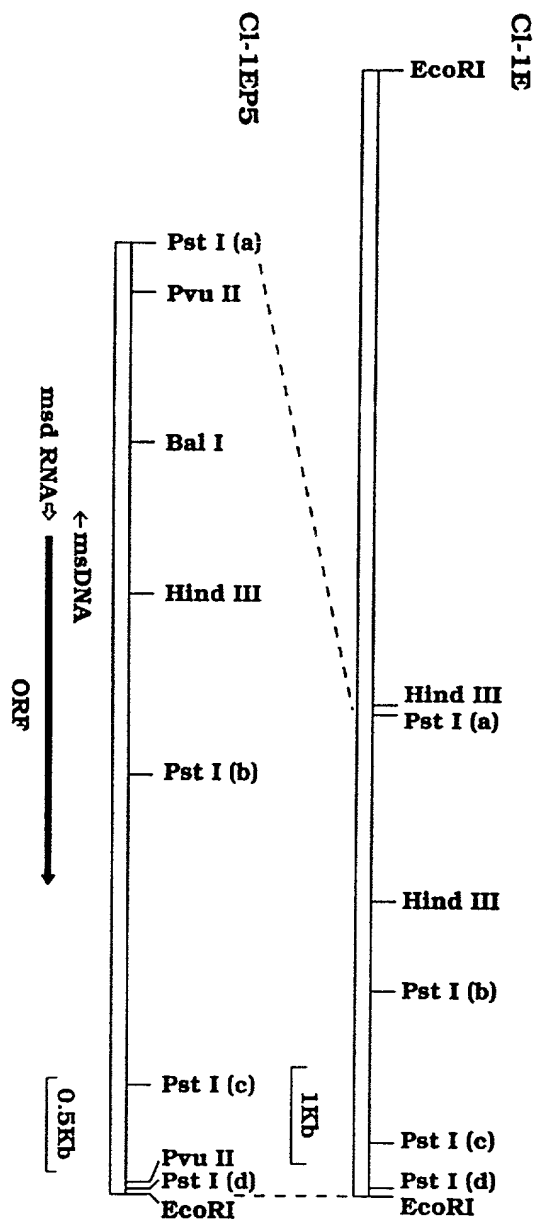


FIGURE 10

08808031.030397

4650ED "TE080830

TCG CGA TT TCA GAT CTT TCA CAG TCG ATG ACT ATG CTG CAT GAA A 120  
CGA TGA TCG ATT GCG GAT CTT TCG TCA GAT CCG CCA GAA CTG CCG CCG TTT TCG TCA 120  
TGT CAT CGA TGT CCA TGA AAA CCA CTG CAT AAA CCG CCG AGG CCG CCG GAT AGC AGC 180  
CGC CCG TAT CAC CGA AAA TAG CCA AAA TAC TTC TCG AAA ACA GAA ACT TGA AGT GAT ATG 240  
TTC ATA AAG AGC GAT CTA CCG ACA TTT GTT GGT TGT GAA TCG GAA CCA CTC CCG TTA ATG 300  
AAG TAT TTG TCG GTA CAT CCG TGT AAA CAA CCA ACA CTT ACC GTT GGT CAC CCG AAT TAG  
CGA CGA GGA ATG CCG TCG CTA AAA TCG TTG ATT CAG AGC TAT ACC GCA GGT CTG CTG TCG 360  
CGT CCG GGT TAG CCG AGC GAT TTT AGC AAG TAA CTC TCG ATA TCG GGT CCA CAG GAG AGC  
GAA GGA CTG CCG CCA TCG GAT TGT CCG TCG CTT TTT TTC CTC TCG CAT GAA GAA GAA ATG 420  
CTT CCG CAC GGA CCG AGC CAA AGA GGA ACC CCA AAA AAG CAG ACC CTA CTT CTT CTT TAC  
ACA AAA ACA TCT AAA CTT GAC GCA CTT AGG CCG CCG ACT TCA CCG CAA GAC TTG CCG AAA 480  
T K T S K L D A L R A A T S R E D L A K  
ATT TTA GAT ATT AAG TTG GTA TTT TTA ACT AAG GTT CTA TAT AGA ATC CCG TCG GAT AAT 540  
I L D I K L V F L T N V L Y R I G S D N  
CAA TAC ACT CAA TTT ACA ATA CCG AAG AAA CCA AAA CCG GTA AGG ACT ATT TGT GCA CCG 600  
Q Y T Q F T I P K K G K G V R T I S A P  
ACA GAC CCG TTG AAG CAC ATC CAA CCA AGA ATA TGT CAC TTA CTT TGT CAT TGT AGA GAT 660  
T D R L K D I Q R R I G D L L S D C R D  
GAG ATC TTT GGT ATA AGC AAA ATT AGT AAG AAC TAT TCG TTT GGT TTT GAG AGC GCA AAA 720  
E I F A I R K I S N N Y S F G F E R C K  
TCA ATA ATC CTA AAT CCG TAT AAG GAT AGA CCG AAA CAA ATA ATA TTA AAT ATA GAT CTT 780  
S I I L N A Y K H R G K Q I I L N I D L  
AAG GAT TTT TTT CAA AGC TTT AAT TTT GGA CCA GTT AGA CCA TAT TTT CTT TCG AAT CAG 840  
K D F F E S F N F G R V R G Y F L S N Q  
GAT TTT TTA TTA AAT CCG CTG GCG CCA AGC AGA CTT CCA AAA CCG CCA TCG TAT AAT CCA 900  
D F L L N P V V A T T L A K A A C Y N G  
ACC CTG CCG CAA CCA AGT CCA TGT TGT CCG ATT ATC TCA AAT CTA ATT TCG AAT ATT ATG 960  
T L P Q G S P G S P I I S N L I C N I H  
GAT ATG AGA TTA CCG AAG CCG GCT AAA AAA TAT CCA TGT ACT TAT AGC AGA TAT CCG GAT 1020  
D H R L A K L A K K Y G C T Y S R A D  
GAT ATA ACA ATT TCT ACA AAT AAA AAT ACA TTT CCG TTA GAA ATG CCG ACT CTG CCA CCG 1080  
D I T I S T N K N T F P L E N A T V Q P  
GAA CCG CTT GTT TTG CCA AAA CTT TTG CTA AAA GAA ATA GAA AAC TGT CCA TTG GAA ATA 1140  
E G V V L G K V L V K E I E N S G F E I  
AAT GAT TCA AAG ACT AGG CTT AGC TAT AAG ACA TCA AGC CAA GAA GAA CCG CCA CTT ACA 1200  
N D S K T R L T Y K T S R Q E V T G L T  
GTT AAC ACA ATC GTT AAT ATT CAT AGA TGT TAT TAT AAA AAA ACT CCG CCG TTG CCA CAT 1260  
V N R I V N I D R C Y Y K K T R A L A H  
CGT TTG TAT CCG ACA GGT GAA TAT AAA CTG CCA GAT GAA AAT CCG GTT TTA GTT TCA GCA 1320  
A L Y R T G E Y K V P D E N G V L V S G  
CGT CTG GAT AAA CTT CAG CCG ATG TTT GGT TTT ATT CAT CAA CTT GAT AAG TTT AAC AAT 1380  
G L D K L E C M F C F I D Q V D K F N N  
ATA AAG AAA AAA CTG AAC AAG CAA CCG GAT AGA TAT GTA TTG ACT AAT CCG ACT TTG CAT 1440  
I K K K L N K Q P D R Y V L T N A T L H  
GGT TTT AAA TTA AAG TTG AAT CCG CCA GAA AAA CCA TAT ACT AAA TTT ATT TAC TAT AAA 1500  
G F K L K L N A R E K A Y S K F I Y Y K  
TTT TTT CAT CCG AAC ACC TGT CCG AGC ATA ATT ACA GAA CCG AAG ACT CAT CCG ATA TAT 1560  
F F H G N T C P T I I T E G K T D R I Y  
TTG AAG CCG CCG TTG CAT TGT TTG CAG ACA TCA TAT CCG GAT TTG TTT GAA GAA AAA ACA 1620  
L K A A L H S L E T S Y P E L F R E K T  
GAT ACT AAA AAG AAA CAA ATA AAT CTT AAT ATA TTT AAA TGT AAT GAA AAG ACC AAA TAT 1680  
D S K K K E I N L N I F K S N E K T K Y  
TTT TTA GAT CTT TGT CCG CCA ACT CCA CAT CTG AAA AAA TTT GTA CAG CCG TAT AAA AAT 1740  
F L D L S G G T A D L K K F V E R Y K N  
AAT TAT GGT TGT TAT TAT CCG TGT CTT CCA AAA CAG CCA GAT ATT ATG GTT CTT GAT AAT 1800  
N Y A S Y Y G S V P K Q P V I M V L D N  
GAT ACA GGT CCA AGC GAT TTA CTT AAT TTT CTG CCG AAT AAA GTT AAA AGC TCG CCA GAC 1860  
D T G P S D L L N F L R N K V K S C P D  
GAT GTA ACT GAA ATG AGA AAG ATA TAT ATT CAT GTT TTC TAT AAT TTA TAT ATA CTT 1920  
D V T E H R K M K Y I H V F Y N L Y I V  
CTG ACA CCA TTG ACT CCG TCG CCG CCA CCA ACT TCA ATG CAG GAT CTT TTC CCG AAA GAT 1980  
L T P L S P S G E Q T S M E D L F P K D  
ATT TTA GAT ATG AAG ATT GAT CCG AAG AAA TTG AAC AAA AAT AAT GAT CCA GAC TCA AAA 2040  
I L D I K I D G K K F N K N N D C D S K  
AGC GAA TAT CCG AAG CAT ATT TTT TCG ATG AGC GTT GTT AGA GAT AAA AAG CCG AAA ATA 2100  
T E Y G K H I F S M R V V R D K K R K I  
GAT TTT AAG CCA TTT TGT TGT ATT TTT GAT GGT ATA AAA GAT ATA AAG CAA GAT TAT AAA 2160  
D F K A F C C I F D A I K D I K E H Y K  
TTA ATG TTA AAT AGC TAA TCA ACA CCG CTA AGC TTA TCA AGC CTA AGC CTG ATT TTT CCG 2220  
L N L N S  
TAA AAT TTA TAT CCG TTG AAT TGT AAT ATA TTA TGT TCA AGC CAT TTA TTT AAT TCG TCG 2280  
ATC GTT TTC TGT AAG CCG ATT AAT TCG TTC CTC ACA AAC ACT AAA CTG CCG TTT TCG ACA 2340  
TCG CCA AAG CCG CCG AAT ATT CCG CAT AAT CCG CAT CAT TTG CCG TCG CAC AGC ATG 2400  
CGC TCG CAT CAT CTC ATG CCG GC

FIGURE 11



V	RT	VKLKPGMDGPKVKQ	WPLTEEKIKALVEICTEMEKEGKISKIGPENPYNTPVFAIKKKDSTKWR	239
LV1	RT	RPWARTPPKAPRNQ	PVPFKPERLQALQHLVRKALEAGHIEPYTG PGNNPVFPVKKA NGTWR	75
DNA	RT	NVLYRIGSDNQYTQFTIPKKGKGVRTISAPTDR	KDIQRRICDLLSDCRDEIFAIRKI SNNYS	94
		+ ○ ● ● ○	+ ●++ ● +	
V	RT	KLVDFRELNKRTQDFWEVQLGIPHPAGLKKK	KSVTVLDVGDAYFSVPLDEDFRKYTAFTIP SI	302
LV1	RT	FIHDLRATNSLTIDLSSSSPGPPDLSSLPTTLAHLQ	TIDLRDAFFQIPLPKQFQPYFAFTVP QQ	139
DNA	RT	FGFE RGKSIILNAYKHRGKQIILNIDLKDDFES	FNFGVRVG YFLS NQDF LLN PVVA	150
		○ ● + ●+ + +○ +● +● ●		
V	RT	NNETPGIRYQYNVLPQGWKGSPAIFQS	SMTKILEPFFKKQNPDIVIYOYMDLLYVGS DLEIG	363
LV1	RT	CNYGPGTRYAWKVLPGGFKNSPTLFEM	QLAHILQPIRQAFPQCTILOYMDILLAS PSHE	199
DNA	RT	TTLAKAACYN GTLPQGSPCSPIISNLICNIMDMRLAKLAKKY	GCTYSRYADDITI STNKNTF	212
		● ●●● ●● + + ○ + ○○ ●●● ●		
V	RT	QHRTKIEELRQHLLRWGLTTP	DKKHQKEP PFLWMGYELHPDKWTVQPIVLPE KDSWTVNDI	424
LV1	RT	DLLLLSEATMASLISHGLPVS	ENKTQQTPTGIKFLGQIISPNHLTYDAVPTVPI RSRWALPEL	262
DNA	RT	PLEMATVQPEGVVLGKVLVKEIENSGFEINDSKTR	LTYSRQEV GTLVNRIVNIDRCYYKKT	276
		○ + ● ○○ ○ + ○ ● + +○		
V	RT	QKLVGKLNWASQIYPGIK	VRQLCKLLRGTKALTEVIPLTEEALELAENREILKEPVGHVYD	487
LV1	RT	QALLGEIQWVSKGTPTLRQPLHSLYCALQRHTDPRD	QIYLNPSQVQSLVQLRQALSQNCRSRLVQ	327
DNA	RT	RALAHALYRTGE YKVPDE NGV	LVSGGLDKLEGMFGFIDQVDKFNNIKKKLNKQ PDRYVL	335
		○● + + + + ○ + ○○ ● ○+○		
V	RT	PSKDLIA EIQKQGQGWTYQIYQE	PFKNLKTGKYARMRGAHTNDVKQLTEAVQKITT	544
LV1	RT	TLPLLGAIMLTITGTTTVVFQSKEQWPLVWLHAPLPHTS	QCPWGQLLASAVLLLDKYTLQSY GL	391
DNA	RT	TNATLHGFKLKL NAREKAY SKFIY YKFFHGNTCPTI	ITEGKTDRIYLKAALHSLET SYPEL	396
		○ ● ○○ + ○○ + ○ ○ + + +○ + ○○ ○		
V	RT	ESIVIWGKTPKFKLPIQKETWETWWTEYWQATWI	PE WEFV NTPPL VKLWYQ	595
LV1	RT	LCQTIHHNISTQTFNQFIQTSDHPSVPILLHSHRFRKNLGA	QTGELWNTFLKTAAPLAPVKALMP	456
DNA	RT	FREKTDSSKKKEINLNIFKSNEKTYFLDLGGTADLKKF	VERYKNNYASYGVS PKQPVIMVLD	460
		+ + ○○ + ○ ○ ●○○		
V	RT	LE KEPIV	GAETFYVDGAANRETKLGKAGYVTNKGRQK VV PLTNTTNQ KTELQAIYLA	652
LV1	RT	VFTLSP VIINTAPCLFSDGSTSRAAYILWDKQILSORS	FP LPPPHKSA Q RAELLGLLHGL	516
DNA	RT	NDTG PSDLLN FLRNKVKSCPDVTEMRKMKYIHVFYNLYIVL	TPLSPSGEQTSMEDLFPKDIL	523
		○ ● ○ + + + ○+○ +○ ●+ ○ ● ○ ● ○		
V	RT	LQDS GLE VNIVTDSQYAL	QIIQA QPDKSESELVNQIIEQLIKKEKVYLAWVPAHKG	708
LV1	RT	SSAR SWR CLNIFLDSKYLYHYLRTLALGTFQGRSSQAPFQA	LLPRLLSRKVVYLHVRSHN	578
DNA	RT	DIKIDGKKFNKNNDGDSKTEYGKHI	FMSR VV RDKKRKIDFKAFCCIFDA	572
		+ ● ●○○ ○ + ○ ○ ○ +○ +		
V	RT	IGGNEQVDKLVSAG		722
LV1	RT	LPDPISRLNALTDA		592
DNA	RT	IKDIKEHYKMLNS		586
		+ ○ ++		

FIGURE 12

**B**

**A**



**M-MuLV**

**pGB2**



**pCI-1EP5**

**FIGURE 13**

00000001-00000001

1

—

2

4

8

6

2

FIGURE 14

FIGURE 15

**RECEIVED**

CC ACT TCC GGC GCT CGG GCT GCG CGA GGG CCC GTG CGA GCA CAT GAT GGC GCT GCG GCT 60  
 GT CCA GGT CGC GCA CGC CGC CGA GGA AGC ACT GCG TCA GAC CCC CGC GGG CGC CCA 120  
 CT CAT CGG CGC GCA GAC GCG CTC CTA GGT GCG GCG CGA CTC GCG CCA GGA GGA GGT 180  
 TA CGG CGT CTC ATT GGA TGG GAA AGT GGT GCG GGT GGA GTG GGG CCC CGC CCA GGG GGA 240  
 TC CGG CGG GCA GAA GCT CTG GTT CGA CAC GGA GCG CGA GGC GCG CAC GCG CTA CTT CAC 300  
 CG OCT GGA GTC CTT GCG GCG GGA GGG ATA TAT CGA TGC GCG TGC TTC AAT GAT GTA GAA 360  
 AC GCA AGC CAC GGG GCG GCG GCG GCG GCG GCG AAA GGC AGG TGC GAC GGA ACG ACA GAC 420  
 22 RNA< 1  
 CT GGT GCG AGC GAC GGA GAG AGG TCC GAA GCG ATC AGC CTC AGC GCG TCG AGC GCG AGC 480  
 CA GCA CGC TCG CTG GCT CTC TCC AGG GTT GCG TAG TCG GAG TCG GCG AGC TCG GCG TCT  
 CG GCG TTG GCG GCG TCT GGT TGA ATT GCA GGA GCG TCT GCG CAA GGT AGC CTG TTC TTG 540  
 CG CCG AAC GCG GCG AGA CCA ACT TAA GGT GGT GAG AGA GCG GTT CCA TCG GAG AAG AAC  
 CT CTC TTC CCG GCG AGT ACC TCT GCG GCG GCG GAG CTG AAC CAA CGA GCG AAC GCG 600  
 CA GAG AAG GGA GCG CAC TCA TGG AGA GCG GCG GCG CTC GAG TAT GGT GGT GCG TTG GCG  
 CT TTC CCG GCG GCG AGA GGT ACT CAC GCG AGG GGA GAG GCG GTG AGG CTA CCG TCC CCG 660  
 CA AAG GCG GCG GCG TCT CCA TGA GTG GCG TCG GGT CTC GCG CAG TCG GAT GCG AGC GCG  
 21  
 CG TGA GAA GGT GGT GCG TTC GGG CCT CCC TCG ACC GGT GCG GCT GCG TCG CCC TCG CCT 720  
 CT ACT CTT CCA CCA CCG AAG CCC GGA GGG AGC TGG CAG GCG CGA GCG AGC GGG AGC GGA  
 CC TCG CCC CCG CCA CCT TCG TCA CCG GCG CCA GGA GCG GTC ATG ACC GCG AAG CTG GAG 780  
 N T A K L E  
 CA CAC GTC CCG GCG GCG CCC CCG GTC TCC GCG GAG GCG CCG GCG CCC ACC GGT CCC GAT 840  
 S H V P A A P P V S A E A P A P T R P D  
 CC GCG AGC CAG GAG GCG GCG GCG CAC CAC GAG GCG CTG GCG CTG GCG TGG AAG GCG 900  
 A A H Q E A R R A H E A L R L R W X A  
 TC GAA GAG GCG GCG GCG ACG GAC GCG TGG GTG GCG CAG CAG CTG GTG GCG AAG GCG GTC 960  
 I E E A G G T D A M V R Q Q L V A K G V  
 CG GCG GCG GAG GTG GAC TTC GAG TCG CTC AGC GAC AAG CAG AAG GCG GCG TGG AAG GAG 1020  
 A A H E V D F E S L S D K Q K A A W K E  
 AG AAG AGC GCG GCG ACC GAG GCG GCG GCG CAG AAG GCG CTG GCG TGG GAG GCG TGG 1080  
 K K H A E A T E R R A Q X R L A W E A W  
 AG GCG AGC CAC ATC CAC CAC CTG GCG GTG GCG GCG CAG TGG GAG GCG GCG GCG CCG 1140  
 K A H I H H L G V G V H W D E A G G P  
 AC AAG TTC GAC GTG GCG GCG GCG GAG GCG GCG AAG GCG AAC GCG TTG CCG GAG GCG 1200  
 D K F D V A G R E E R A X A N G L P E G  
 TG GAC TCG GTC GAG GCG CTG GCG AAA GCG CTG GCG ATC TCC GTG TCG GCG CTG GCG TGG 1260  
 L D H V E A L A K A L G I S V S R L R W  
 TC TCC TCG CAC GCG GAG GTG GAC AGC GCG AGC CAC TAC CAG AGC TGG GAG ATT CCG AAG 1320  
 F S H R E V D T G T H Y Q T W E I P K  
 CG GAG GCG GCG AAG GCG ACG CTC ACC GCG GCG AAG GCG CAG CTC AAG GCG CTG CAG CCG 1380  
 R D G G K R T L T A P K R E L K A V Q R  
 CG GTG CTC GCG AAC GTG GAG GCG CTG CCG GTG CAC GCG GCG GCG CAC GCG TTC GTG 1440  
 W V E A N V V E R L P V H G A A H G F V

GCG GCG GCG TCC ATC CTC ACC AAC GCG CTG GCG CAC CAG GCG GCG GAC GTG GTG GTG AAG 1500  
 A G R S I L T M A L A H Q G A D V V V K  
 GTG GAC ATG AAG GAC TTC TTC OCT TCC GTG ACG TGG CCC CGG GTC AAG GGA CTG CTG CCG 1560  
 V D M K D F P S V T W P R V K G L L R  
 250  
 AAG GGA GGA CTC CCG GAG AAC CTC GCG ACG CTC CTG GCG CTG CTC TCC ACC GAG GCG CCG 1620  
 K G G L P E N L A T L L A L L S T E A P  
 CGC GAG GTG GTG CCG TTC CCG GGA GAG ACG CTG TAC GTG GCG AAG GCG OCT GCG GCG CTG 1680  
 R E V V R F R G E T L Y V A K G P R A L  
 300  
 CCC CAG GCG GCG CCC ACC TCT CCG GCG CTG ACG AAC GCG CTG TCG CTG GCG CTG GAC AAG 1740  
 P Q G A P T S P A L T N A L C L R L D K  
 420  
 CCG CTC TCG GCG CTG TCG AAG CCG CTG GCG TTC ACG TAC ACG GCG TAT CCG GAT GAC CTG 1800  
 R L S A L S K R L G F T Y T R Y A D B L  
 540  
 ACG TTC TCC TCG CCG CCG GCG AAG AAG TCC CCG CAG AAG GAA CTC CCC CTG GCG GAT GCG 1860  
 T F S W R R A K X S R Q K E L P L A D A  
 600  
 CCG GTG GCG CTG CTC CTG GCG GCG GTG AAG GGT GTG CTG GAG GCG GAG GGT TTC ACG CTG 1920  
 P V A L L L A R V K G V L E A E G F T L  
 660  
 CAC CCG GAC AAG ACG CCG GTG CAG CCG AAG GCG AGC CCG CAG CCG GTG ACG GCG CTC GTG 1980  
 H P D K T R V Q R K G S R Q R V T G L V  
 720  
 GTG AAC GAG GCG CCC GAG GCG GTT CCG GGT GCG CCG GTG CCC CCG GAT GTG GTG CCG GCG 2040  
 V N E A P E G V P G A R V P R D V V R R  
 780  
 CTG CCG GCG GCG ATC CAC AAC CCG GAG CAG GCG AAG CCC GCG CCC ACC GCG GAG ACG CTG 2100  
 L R A A I H N R E Q G K P G P T G E T L  
 840  
 GAG CAG CTC AAG GCG CTC GCG GCG TTC CTT CAC ATG ACG GAC GCG GAG AAG GCG GCG GCG 2160  
 E Q L K G L A A F L H M T D A E K G R A  
 900  
 TTC CTG CGA CCG CTG GAG GCG CTC GAG AAG CCG CAG ACC GCG TGA CCC TCA CTG GTC GTC 2220  
 F L R R L E A L E K R Q T A -  
 960  
 CCG GCG ATC GCA GCG GCG GCG ACG GAC COT CAC CCC CCA GAT CTC CAT GCG ATG CTG 2280  
 GGG ATT CTG GCG GGT GAA GAA GAC TTC CCA GCG GAG ACG GAC GAA GCG CTG GCG ATC CGA 2340  
 TCA CTC CTC GCG GCG GAT CTC CCG GAG GCG CAG GGT TCC GAC GTC COT GCG ATT GCT 2400  
 1020  
 CAC CCA GCG CTC CCG GCG CCA GCG TTG GGT GTC CCG CGA GAA GAA GAG CAG CCC GGA GAT 2460  
 GCG COT CAG GTT CTC CCG CGA GCG ATC CTC GCG GCG GCG CAA ATC CTT CAG CAG CAG 2520  
 1080  
 GGT GCG CTT GCG GGT GCG ATC GCT GGA CCA CAG CTC CCG GCG GTG GAG GGT GTC ACT CCG 2580  
 GCG GAA GTA GAG CAT CCC ATT CAG CCG CTT GAT GCG GGT GCG GCG CCA GCT GTC CCG ACC 2640  
 1200  
 CCG CCA GAT GTC CTT CAC CCG GAC GGT GCG ATG CGA GGT GCG ATC GCT GAC CCA CAG CTC 2700  
 CTC GCG CTC GCG CTG GCG CCA GAA CTC GCG CTC GCG TCC CCG GCG GCT GAA GAA GAT CTT 2760  
 1260  
 CCC CCC GAG CCG COT GAG ATC ATG CCG ATA GAG GCG GCG GAA GAA GCG CAG CTG CTC GGA 2820  
 GAC GGT GCG TCT GGA GCA CCA CAG GCT GCG CTC GCG TTC GTC ATT GTC GAG CAG GAA GAA 2880  
 1320  
 GAG CAC CGA GTC CCG CCG GGT GAA CCG GGA GAG GAA GTT GTC CTC GCG GCG COT GAA GAC 2940  
 AGA COT GGT GCT GGA CAG CCG CAG GGT GCG CCA GAT GAA CAC CTC GTC ATT CAC GTT GCG 3000  
 1380  
 CAC GAA GAA GAG CCG ATC GCG GAC CCG GGT GAG CCG GCG GCG GCT GGA GCT GCG GCG CAC 3060

FIGURE 16

030397

TTTCGAGAAG CGCCATACCA AACAGGGGAT ACAGACCAAC CTGACGCTGA AAGAGGAAG CTACGGGAC TGCGTCCGA AGTCCGAAGA 9990  
F E K R N T K Q G I Q T N L T L K E E S Y G D M L P K C D D  
CCCCGACGA ACATAACCTC ACTCAGACCG GCAACAGCGG GTCTTTTCTT TTCTGGCCAT TCCGACAAGG TGAACAATCC ACTGTTTACC 10080  
P A A T \*  
CTTCACCGTT TATTCACCTT TTATCACTAT GAAATTATTA ATAAAAAACC AGAGGTGAC ACTGTCACCA GTAAACCTG AAAAACTTT 10170  
TTATCACCCC CGGCATCGCC CGACTCGACA GATCCAGAAC GAGCAAAAT CACAAAGGTG ACCAGTCCAC TGTCTACTCT TCACCAACTC 10260  
ATCACCACCT AACCAATGTA TATAAATGA TAAATAATCG AGGTGACAG TTAATGCGA AAAAAGCTTT TCTCAGCTCT TGGATAAAG 10350  
AAAAATTAAT CACATCAATA GCTTCTCTCT TGAATCTCTT TGAGGTTTAT GAGAGCGTAA TAGAGCCAAA GCTTCAATTT TATGGGTTAT 10440  
TTTTAATTAAT GTGTAGTTAT CGAAAGGAGA ACTTAGGAGA ACTCCAAATA CTCTGCGATT GTCTCGGTTT GGATCGTAAA ATACCCAATT  
TAGCCCATCG CGCATGAGTC ATGTTTTCGC CTAGTATTTT AGCTATGCCC GTGCTTCAGT TCGCTGAGCG CGCGCTGGGG GGCACCGATC 10530  
ATCGGCTAGC GGTACTCTAG TACCAGAGCG CATCATAAAA TCGATACGGG CAGCAAGTCA AGCGACTGCG CGCGACCCCG CGCTGCTAG  
AGCGAATCGA TCGAGCTGCT CAGTAGGTTT TGGCTCTTTT AGTCTCTTAC CATCAAGGTG CATAAGGATA TTCTCGATGC TGACTCAGCT 10620  
TCGCTTGACT AGCTGACGCA GTTATCCAA ACCGAGAAAA TCAGGAGATG GTAGTTCCAC GTATTCTTAT AAGAGCTAGC ACTGATCCCA  
or1316 M L T Q L  
AAAAAAAAT GGTACTGAGG TATCTAGAGC AACCGGCTTA TTTTCATCAT TCGTTGAAAG GAACAAAGTA AAATGCTCTG GTAAATGAAA 10710  
K K N G T E V S R A T A L F E S F V E K N X V K C P G N V K  
AAAAATGCTT TTCTGTGTG GTGCTAACAA AAACAATGGA GAACCATCAG CAAGCAGATT GGAATTAATA AATTTTCTG AAAGGTATT 10800  
K F V F L C G A N K N M G E P S A R R L E L I N F S E R L  
GAATAAATCT CACTTTTTC TTGCTGAAT AGTTTTCAAA GAATTAAGCA CCGATGAAGA ATCATTATCT GATAATTTAT TAGATATCGA 10890  
N N C H F F L A E L V F K E L S T D E E S L S D N L L D I E  
AGCTGACTTA TCTAAATTAG CTGATCATAT TATCATGTTT TTAGAAAGTT ATTCACTCTT CACGGAACTT GGTGCAATCG CATACAGCAA 10980  
A D L S K L A D N I I I V L E S Y S S F T E L G A F A Y S K  
GCAATTCGC AAGAAATTA TAATAGTTAA CAATACAAAA TTTATAATG AGAAATCATT TATAAATATG CGACCAATAA AGGCTATTAC 11070  
Q L R K K L I I V N N T K F I N E X S F I N N G P I K A I T  
TCAGCAATCA CAACAATCTG GTCATTCTTT ACATTATAA ATGACAGAAG GTATTGAAG TATAGAGCGC TCTGATGGA TTGCGGAAAT 11160  
Q Q S Q S G H F L H Y K N T E G I E S I E R S D G I G E I  
ATTCGACCCC CTATATGATA TTCTTTCTAA GAAGCAGAGA GCAATTTCAA GAACTTAAAA AAAAGAGAGG TTAGATCTCT CCAGTAACCT 11250  
F D P L Y D I L S K N D R A I S R T L K K E E L D P S S M F  
CAATAAGAC TCAGTACGAT TTATTCATGA CGTAATTTT GTATGTGCTC CTTTGCAACT TAATGAATC ATCGAAATTA TCACAAATAT 11340  
N K D S V R F I N D V I F V C G P L Q L N E L I E I I T K I  
ATTGCGACA GAAAGCCATT ACAAAAAAAA TCTCTAAGG CACTTTGGTA TTCTAATAGC TATTAGAATA ATATCATGCA CAATGGGAT 11430  
F G T E S H Y K K N L L K N L G I L I A I R I I S C T N G I  
TTATTATCT TTGTATAAG AATATTATTT TAAATATGAC TTGACATTG ACAAATATC ATCAATGTTT AAAGTTTITT TCTCAAGAA 11520  
Y Y S L Y K E Y Y F K Y D F D I D M I S N F K V F F L K N  
CAAGCCAGAA AGGATGAGGG TATATGAGAA TATATGACCT AATTGATTCT CAGACATTGA TGACTAAGGG ATTGCTTCTT GAAGTAATGC 11610  
K P E R N R V Y E N I \*  
RT N R I Y S L I D S Q T L N T K G F A S E V N  
GATCACCTGA GCGGCAAAA AAATGGGATA TAGTAAGAA AAAAGGAGGT ATGAGAACAA TTATCACCCC GTATCAAAA GTTAAATTA 11700  
R S P E P P K K M D I A E K K G G N R T I Y N P S S K V K L  
TTCAATATGG GTTAATGAAT AATGTTTTTT CGAAGCTCCC AATGCATAAT GCTGCATATG CATTGTGTTA AAACCGATCA ATAAAAAGCA 11790  
I Q Y W L M K N V F S K L P M H N A A Y A F V K N R S I K S  
ATGCTTATT ACATGGCGAA TCAAGAATA AGTATTATCT GAAATAGAT CTCAAGATT TTTCCTCTT AATAAATTT ACTGATTGTT 11880  
N A L L N A E S K N K Y Y V K I D L K D F F P S I X F T D F  
AGTACGATT CACTGTTAT CGAGTCGCA TGAATTTAC TACAGAAAT GATAAGGAGT TACTACAAT TATAAAGG ATCTGCTTTA 11970  
E Y A F T R Y R D R I E F T T E Y D K E L L Q L I X T I C F  
TATCAGATAG CACTCTCCCT ATCGGGTTTC CTACATCTCC ATTAATTGCA AACTTTGTCG CAAGAGAACT TGATGAAAA CTGACGCAAA 12060  
I S D S T L P I G F P T S P L I A N F V A R E L D E K L T Q  
AACTAAATGC AATTGATAAA CTTAATGCA CTATACAGC ATATGCTGAT GATATTATG TCTCTACAAA TATGAAGGG GCTAGCAAT 12150  
K L M A I D K L M A T Y T R Y A D D I I V S T N M K G A S E  
TAATCTGGA TTGTTTTTAA AGAACAATGA AAGAGATTGG TCCAGACTTT AAAAATTAACA TTAATAAAT TAAGATTGTT AGTCTCTCG 12240  
L I L D C P K R T M K E I G P D F K I M I E K F K I C S A S  
GAGGAATAT AGTAGTTACC GGATTGAAG TTTGCCAGA TTTTCATATT ACATTACATA GATCAATGAA AGATAAATA AGATTGCATC 12330  
G G S I V V T G L K V C N D F H I T L H R S M R D X I R L N  
TTTCTCTTT ATCAAAGGCG ATATTAAAG ATGAAGATCA TAATAAACTT TCTGTTTATA TTGCTTATGC AAAAGATATA GACCTCAT 12420  
L S L L S K G I L K D E D M N K L S G Y I A Y A K D I D P N  
TTTATACAA ACTGAACAGA AATATTATTC AAGAAATAA ATGGATTACG AATCTCCACA ACAAAGTTGA ATAACTTTA TATTTGGAT 12510  
F Y T K L M R K Y F Q E I K M I Q N L N N K V E \*  
GCACCCCAAT AACTTCATTG ATTAATGCG GAACAATATA GGTCTTTCAG GATGACCTAC ACTCTAGAGA ATGTGTATAC AAAAGTGTAT 12600  
AAGTTATTTT CAAACCTATA TAAATAACAG CAAATCAAT GCATTGGCGG CATTTTACCG CTCTGTGAT CTTCGCGCAA AATGCTG 12688

FIGURE 17

Year	1970	1971	1972	1973	1974	1975	1976	1977	1978	1979	1980	1981	1982	1983	1984	1985	1986	1987	1988	1989	1990	1991	1992	1993	1994	1995	1996	1997	1998	1999	2000	2001	2002	2003	2004	2005	2006	2007	2008	2009	2010	2011	2012	2013	2014	2015	2016	2017	2018	2019	2020	2021	2022	2023	2024	2025	2026	2027	2028	2029	2030	2031	2032	2033	2034	2035	2036	2037	2038	2039	2040	2041	2042	2043	2044	2045	2046	2047	2048	2049	2050	2051	2052	2053	2054	2055	2056	2057	2058	2059	2060	2061	2062	2063	2064	2065	2066	2067	2068	2069	2070	2071	2072	2073	2074	2075	2076	2077	2078	2079	2080	2081	2082	2083	2084	2085	2086	2087	2088	2089	2090	2091	2092	2093	2094	2095	2096	2097	2098	2099	2100
1970	1971	1972	1973	1974	1975	1976	1977	1978	1979	1980	1981	1982	1983	1984	1985	1986	1987	1988	1989	1990	1991	1992	1993	1994	1995	1996	1997	1998	1999	2000	2001	2002	2003	2004	2005	2006	2007	2008	2009	2010	2011	2012	2013	2014	2015	2016	2017	2018	2019	2020	2021	2022	2023	2024	2025	2026	2027	2028	2029	2030	2031	2032	2033	2034	2035	2036	2037	2038	2039	2040	2041	2042	2043	2044	2045	2046	2047	2048	2049	2050	2051	2052	2053	2054	2055	2056	2057	2058	2059	2060	2061	2062	2063	2064	2065	2066	2067	2068	2069	2070	2071	2072	2073	2074	2075	2076	2077	2078	2079	2080	2081	2082	2083	2084	2085	2086	2087	2088	2089	2090	2091	2092	2093	2094	2095	2096	2097	2098	2099	2100	

FIGURE 18

0308031-030397

Oligo 2337  
tcaccctgaaagacctgattgcttacctggaagagaagccggaatggcggaacatctgg 60  
cggcgggttaaggccctatcggaagagttcggcgtttataaAATATGCGCTGTGCAGGGTTT 120  
TTGCTGTGCGCAGCGTGATGCGCTTCAAGATATCGTGTAAATCTGCTTTGCGCAGCATG 180  
AACGACACGGCTCGCACTACGCGAAGTTCTATAGCACAATTAGACGAAAGCGGTCTCAC  
GCAATACGTTTCCGGCCTTTTGTGCCGGGAGGGTCGGCGAGTCTGCTGACTTAACGCCAG 240  
CGTTATCGCAAAGGCCGAAAACACGGCCCTCCAGCCGCTCAGCGACTGAATTGCGGTC  
TAGTATGTCCATATACCCAAAGTCGCTTCATTGTACCTGAGTACGCTTCGCGTACGTGCG 300  
ATCATAAGGTATATGGGTTTCAGCGAAGTAACATGGACTCATGCGAAGCGCATGACGCG  
GCTGACGCGCTCAGTACAGTTACGCGCCTTCGGGATGGTTTAATGGTATTGCCGCTGTTG 360  
CGACTGCGCGAGTCATGTCAATGCGCGAAGCCCTACCAATTAACATAACGCCGACAAC  
GGCCCTCTTTTGGCCGCGTGATGTGGAGAGTGAATGGATGCTACCCGGACAACCCCTTC 420  
M D A T R T T L L  
TGGCGCTCGATTGTTCGGCTCGCCGGGCTGGAGCGCCGATAAAGAAATACAGCGACTGC 480  
A L D L F G S P G W S A D K E I Q R L H  
ATGCGCTCAGTAATCATGCCGACGCCATTACCGACGCATTATTCTTTCTAAAGCCACG 540  
A L S N H A G R H Y R R I I L S K R H G  
GTGGTCAGCGGCTGGTGTAGCCCTGATTACTTGCTCAAAACCGTACAGCGCAACATTC 600  
G Q R L V L A P D Y L L K T V Q R N I L  
TTAAGAACGTCCTTTTCACAATTTCCGCTTTCCCTTTTGTACAGCCTACCGACAGGTT 660  
K N V L S Q F P L S P F A T A Y R P G C  
GCCCAATCGTCAGCAACCGCGACGCCACTGCCAACAGCCGAGATCCTGAAACTCGATA 720  
P I V S N A Q P H C Q Q P Q I L K L D I  
TCGAAAACCTTTTCGATAGCATTAGCTGGTTACAGGTCTGGCGTGTGTTTCGCCAGGCCC 780  
E N F F D S I S W L Q V W R V F R Q A Q  
AGTTGCCACGTAATGTGGTAACCATGCTGACCTGGATTGTTGTTATAACGACGCGTTAC 840  
L P R N V V T M L T W I C C Y N D A L P  
CGCAGGGGGCACCAACTTCGCCAGCCATTCCAATCTTGTGATGCGCGTTTTGATGAAC 900  
Q G A P T S P A I S N L V M R R F D E R  
GCATAGGGGAATGGTGTGAGGCTCGGGGAATTACCTACACCCGCTACTGCGATGACATGA 960  
I G E W C Q A R G I T Y T R Y C D D M T  
CCTTTTCAGGTCACTTCAATGCCCGCCAGGTTAAAAATAAAGTGTGCGGATTGTTAGCGG 1020  
F S G H F N A R Q V K N K V C G L L A E  
AGCTGGGCTGAGCCTCAATAAACGCAAGGCTGCCTGATAGCTGCCTGTAAGGCCAGC 1080  
L G L S L N K R K G C L I A A C K R Q Q  
AAGTAACCGGGATTGTTGTTAATCACAAGCCACAGCTTGCCCGTGAAGCGCGCGGGCGC 1140  
V T G I V V N H K P Q L A R E A R R A L  
TGCGTCAGGAGGTGCATTGTGCCCCAAATATGGCGTTATTTGCGATCTTAGTCATCGTG 1200  
R Q E V H L C Q K Y G V I S H L S H R G  
GTGAACCTGATCCTTCTGCGATCTCCACGCACAGGCAACGGGTATCTTTATGCTTTGC 1260  
E L D P S G D L H A Q A T A Y L Y A L Q  
AGGGAAGAATAAACTGGTTATTGCAATCAACCCTGAGGATGAGGCCTTTCAACAGGCGA 1320  
G R I N W L L Q I N P E D E A F Q Q A R  
GAGAGAGTGTAAAGCGAATGCTGGTTGCATGGTAAGAAAAGCGTCAGGCAGACGTTTCTG 1380  
E S V K R M L V A W \*  
CCTGACCGTTTAGGGGAGAattactgcaactgcgcggaattagcgccagcgggcggtca 1440  
aatcatccgtcgggcggtattttaaactcgctgcggaacaaacgtgacagcataccttca 1500  
cagaagccaggatctggttgcagcaggggttcatcgg 1540  
Oligo 2336

FIGURE 19



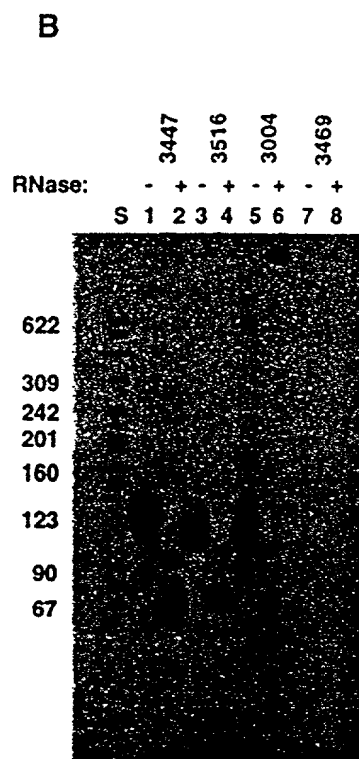
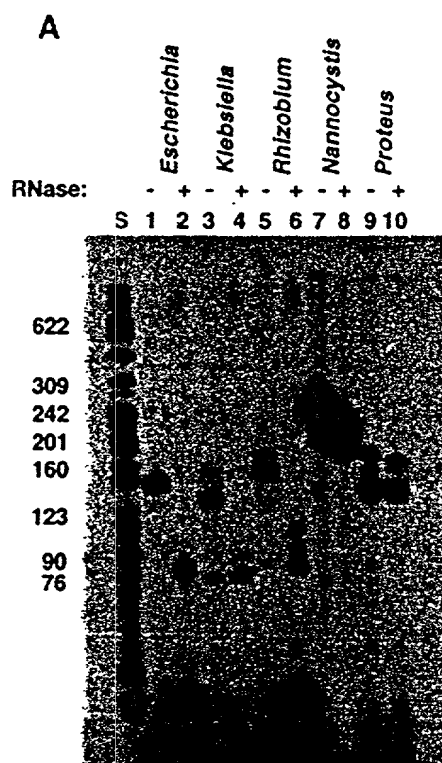


FIGURE 20

## RHIZOBIAL ISOLATES

Strain (legume host genus)	USDA strain no.	Geographic source (date)	msDNA produced <sup>b</sup>
<i>Rhizobium</i> sp. ( <i>Acacia</i> )	3002	Brazil (1959)	+
	3003	Africa (1950)	
	3325	Morocco (1974)	
	3838	? (1976)	+
<i>Bradyrhizobium</i> sp. ( <i>Aeschynomene</i> )	3516	Florida (1972)	+
	4362		
<i>Bradyrhizobium</i> sp. ( <i>Albizia</i> )	3004	Maryland (1952)	+
<i>Bradyrhizobium</i> sp. ( <i>Apios</i> )	3240	Maryland (1939)	
<i>Bradyrhizobium</i> sp. ( <i>Arachis</i> )	3339	Thailand (1979)	
	3341	Hawaii (1978)	
<i>Rhizobium</i> sp. ( <i>Astragalus</i> )	3854	Alaska (1962)	
<i>Rhizobium</i> sp. ( <i>Cajanus</i> )	3472		
<i>Bradyrhizobium</i> sp. ( <i>Canavalia</i> )	3317	Brazil (1974)	
<i>Rhizobium</i> sp. ( <i>Cicer</i> )	3378		
	3379	Mexico (1963)	
<i>Bradyrhizobium</i> sp. ( <i>Coronilla</i> )	3165	Virginia (1935)	
	3167	? (1961)	
<i>Bradyrhizobium</i> sp. ( <i>Crotalaria</i> )	3384	Brazil (1967)	
<i>Bradyrhizobium</i> sp. ( <i>Desmodium</i> )	3225	Ecuador (1948)	
<i>Bradyrhizobium</i> sp. ( <i>Erythrina</i> )	3241		
	3242	Maryland (1939)	+
<i>Rhizobium fredii</i>	191	China (1979)	
<i>Rhizobium leguminosarum</i>	2370	Illinois (1933)	
	2429	Hawaii (1978)	
	2435	Holland (1955)	
	2480	Tennessee (1951)	
	2489		
<i>Rhizobium</i> sp. ( <i>Lens</i> )	2426		
	3404	Colombia (1979)	
<i>Rhizobium loti</i>	3084	Maryland (1946)	
	3468	New Zealand (1961)	+
	3469		
	3471		
	3503		+
	3669	California (1968)	
<i>Bradyrhizobium</i> sp. ( <i>Lotus</i> )	3074	Minnesota (1954)	
	3470	California (1916)	
<i>Rhizobium</i> sp. ( <i>Lupinas</i> )	3040	Florida (1940)	
<i>Bradyrhizobium</i> sp. ( <i>Lupinas</i> )	3045	Florida (1946)	
<i>Bradyrhizobium</i> sp. ( <i>Macrotyloma</i> )	3451	Zimbabwe (1960)	
<i>Rhizobium medicago</i>	1097	North Dakota (1948)	
<i>Rhizobium meliloti</i>	1011	Maryland (1933)	
	1021a	North Dakota (1948)	
<i>Rhizobium phaseoli</i>	2667	Washington (1948)	
	2669		
	2674	Brazil (?)	
	2676	Colombia (1972)	
	3256	Illinois (1941)	
<i>Rhizobium</i> sp. ( <i>Robinia</i> )	3436		
<i>Bradyrhizobium</i> sp. ( <i>Stylosanthes</i> )	3441	Brazil (?)	
	3477	Colombia (1976)	
<i>Rhizobium trifolii</i>	2046	Virginia (1934)	
	2048	Illinois (1934)	+
	2063	Florida (1939)	
	2065	Alabama (1952)	+
	2116	South Carolina (1944)	
	2134	? (1974)	
	2145		
	2156	California (1920)	
<i>Rhizobium</i> sp. ( <i>Trigonella</i> )	1177	Florida (1939)	
<i>Rhizobium tropici</i>	2744	Brazil (?)	
<i>Bradyrhizobium</i> sp. ( <i>Vigna</i> )	3447	Thailand (1979)	+
	3456	Wisconsin (1966)	

<sup>a</sup> All strains are from the USDA Beltsville Rhizobium Culture Collection, provided by Peter van Berkum.

<sup>b</sup> As defined by detection of radiolabeled msDNA by the RT extension method.

SEQUENCE LISTING

(1) GENERAL INFORMATION:

(i) APPLICANT: Inouye, Sumiko  
Hsu, Mei-Yin  
Eagle, Susan  
Inouye, Masayori

(ii) TITLE OF INVENTION: Prokaryotic Reverse Transcriptase

(iii) NUMBER OF SEQUENCES: 45

(iv) CORRESPONDENCE ADDRESS:

(A) ADDRESSEE: Weiser & Associates  
(B) STREET: 230 South Fifteenth Street, Suite 500  
(C) CITY: Philadelphia  
(D) STATE: Pennsylvania  
(E) COUNTRY: U.S.A.  
(F) ZIP: 19102

(v) COMPUTER READABLE FORM:

(A) MEDIUM TYPE: Floppy disk  
(B) COMPUTER: IBM PC compatible  
(C) OPERATING SYSTEM: PC-DOS/MS-DOS  
(D) SOFTWARE: PatentIn Release #1.0, Version #1.25

(vi) CURRENT APPLICATION DATA:

(A) APPLICATION NUMBER: US 08/269,118  
(B) FILING DATE: 30-JUN-1994  
(C) CLASSIFICATION:

(viii) ATTORNEY/AGENT INFORMATION:

(A) NAME: Weiser, Gerard J.  
(B) REGISTRATION NUMBER: 19,763  
(C) REFERENCE/DOCKET NUMBER: 377.5888P

(ix) TELECOMMUNICATION INFORMATION:

(A) TELEPHONE: 215-875-8383  
(B) TELEFAX: 215-875-8394

(2) INFORMATION FOR SEQ ID NO:1:

(i) SEQUENCE CHARACTERISTICS:

(A) LENGTH: 2176 base pairs  
(B) TYPE: nucleic acid  
(C) STRANDEDNESS: double  
(D) TOPOLOGY: linear

(ix) FEATURE:

(A) NAME/KEY: CDS

2630601-030397

(B) LOCATION: 640..2094

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:1:

TCATCCGCGC GGACACCCCC TCCTACGTGC CCCCCGACGC GGAGAGCGGC GTGGAGACGG 60  
TGTACCGCGT TTCCCTGGAT GGTACACCTGG TGGCGGTGGA GTGGGGCCCCG CGCACGGGCT 120  
CGCCGCGTCA CCAGCGGCTC TGGTTCGACT CGGATGCGGA AGCCCCCGGA GCCTACTTCG 180  
CGCGCCTCGA GAAGTTGGCG GCTGACGGCT ACATCGACGC GGCCTCGGCA TTGGTCTAAA 240  
CCCTTCAACC ACGGCTCGGC CGCCACGCGC GGCCGGCAGG ACAGGTGCGA CGAACAGACG 300  
ACGACGTGCG CTTCACGCGC GAGCAGCCGA GAGAGGTCCG GAGTGCATCA GCCTGAGCGC 360  
CTCGAGCGGC GGAGCGGCGT TCGCCGCTC CGGTTGGAAT GCAGGACACT CTCCGCAAGG 420  
TAGCCTGTTC TTGGCTCTCT CCCTCCTAGG CACTACGGCC AGGGTGGGTA GCGGAGCCAA 480  
CGACGCCACC GCCGTTTACC CACCCCGGCC GTAGTGCCTA GGAGGGGAGA GCCGGTGAGG 540  
CTACCGTGCC CCAGGTAAGA TGGTGGTGCT TTCCCGGCCT CCGTCGACTG CTCGCGCCAT 600  
GTCCCGTCTT CCATCGCCGC GCCCGCCCAA GGTGCAGAC ATG ACC GCC AGG CTG 654  
Met Thr Ala Arg Leu  
1 5  
GAC CCG TTC GTC CCC GCA GCT TCG CCG CAG GCC GTG CCC ACG CCC GAG 702  
Asp Pro Phe Val Pro Ala Ala Ser Pro Gln Ala Val Pro Thr Pro Glu  
10 15 20  
CTC ACC GCT CCG TCG TCA GAC GCG GCC GCG AAG CGT GAA GCC CGC CGG 750  
Leu Thr Ala Pro Ser Ser Asp Ala Ala Lys Arg Glu Ala Arg Arg  
25 30 35  
CTC GCG CAC GAA GCG TTG CTC GTC CGC GCG AAG GCC ATC GAC GAA GCG 798  
Leu Ala His Glu Ala Leu Leu Val Arg Ala Lys Ala Ile Asp Glu Ala  
40 45 50  
GGC GGC GCC GAC GAC TGG GTG CAG GCG CAG CTC GTC TCC AAG GGG CTC 846  
Gly Gly Ala Asp Asp Trp Val Gln Ala Gln Leu Val Ser Lys Gly Leu  
55 60 65  
GCG GTC GAG GAC CTG GAC TTC TCC AGC GCC TCC GAG AAG GAC AAG AAG 894  
Ala Val Glu Asp Leu Asp Phe Ser Ser Ala Ser Glu Lys Asp Lys Lys  
70 75 80 85  
GCC TGG AAG GAG AAG AAG AAG GCC GAG GCC ACC GAG CGC CGC GCG CTG 942  
Ala Trp Lys Glu Lys Lys Lys Ala Glu Ala Thr Glu Arg Arg Ala Leu  
90 95 100  
AAG CGT CAG GCG CAC GAG GCG TGG AAG GCC ACG CAC GTG GGC CAC CTG 990

Lys	Arg	Gln	Ala	His	Glu	Ala	Trp	Lys	Ala	Thr	His	Val	Gly	His	Leu	
			105					110					115			
GGC	GCG	GCC	GTG	CAC	TGG	GCG	GAG	GAC	CGC	CTG	GCC	GAC	GCG	TTC	GAC	1038
Gly	Ala	Gly	Val	His	Trp	Ala	Glu	Asp	Arg	Leu	Ala	Asp	Ala	Phe	Asp	
		120					125					130				
GTG	CCC	CAC	CGC	GAG	GAG	CGC	GCC	CGG	GCC	AAC	GGC	CTG	ACG	GAG	CTG	1086
Val	Pro	His	Arg	Glu	Glu	Arg	Ala	Arg	Ala	Asn	Gly	Leu	Thr	Glu	Leu	
	135					140					145					
GAC	TCC	GCG	GAG	GCG	CTG	GCC	AAG	GCG	CTG	GGG	CTG	AGC	GTC	TCC	AAG	1134
Asp	Ser	Ala	Glu	Ala	Leu	Ala	Lys	Ala	Leu	Gly	Leu	Ser	Val	Ser	Lys	
150					155					160					165	
CTC	CGC	TGG	TTC	GCG	TTC	CAC	CGG	GAG	GTC	GAC	ACG	GCC	ACG	CAC	TAC	1182
Leu	Arg	Trp	Phe	Ala	Phe	His	Arg	Glu	Val	Asp	Thr	Ala	Thr	His	Tyr	
			170					175						180		
GTG	AGC	TGG	ACC	ATT	CCG	AAG	CGG	GAC	GGC	AGC	AAG	CGC	ACG	ATT	ACG	1230
Val	Ser	Trp	Thr	Ile	Pro	Lys	Arg	Asp	Gly	Ser	Lys	Arg	Thr	Ile	Thr	
			185					190					195			
TCC	CCC	AAG	CCT	GAG	CTG	AAG	GCA	GCG	CAG	CGC	TGG	GTG	CTG	TCC	AAC	1278
Ser	Pro	Lys	Pro	Glu	Leu	Lys	Ala	Ala	Gln	Arg	Trp	Val	Leu	Ser	Asn	
		200					205					210				
GTC	GTG	GAG	CGG	CTG	CCG	GTC	CAC	GGC	GCC	GCC	CAC	GGC	TTC	GTG	GCG	1326
Val	Val	Glu	Arg	Leu	Pro	Val	His	Gly	Ala	Ala	His	Gly	Phe	Val	Ala	
	215					220					225					
GGA	CGC	TCC	ATC	CTC	ACC	AAC	GCG	CTG	GCC	CAC	CAG	GGC	GCG	GAC	GTC	1374
Gly	Arg	Ser	Ile	Leu	Thr	Asn	Ala	Leu	Ala	His	Gln	Gly	Ala	Asp	Val	
230					235					240					245	
GTG	GTC	AAG	GTG	GAC	CTC	AAG	GAC	TTC	TTC	CCC	TCC	GTC	ACC	TGG	CGC	1422
Val	Val	Lys	Val	Asp	Leu	Lys	Asp	Phe	Phe	Pro	Ser	Val	Thr	Trp	Arg	
				250				255						260		
CGG	GTG	AAG	GGC	CTG	TTG	CGC	AAG	GGC	GGC	CTG	CGG	CAG	GGC	ACG	TCC	1470
Arg	Val	Lys	Gly	Leu	Leu	Arg	Lys	Gly	Gly	Leu	Arg	Glu	Gly	Thr	Ser	
			265					270				275				
ACG	CTG	CTG	TCC	CTC	CTC	TCC	ACG	GAA	GCG	CCG	CGG	GAG	GCG	GTC	CAG	1518
Thr	Leu	Leu	Ser	Leu	Leu	Ser	Thr	Glu	Ala	Pro	Arg	Glu	Ala	Val	Gln	
		280					285					290				
TTC	CGC	GGC	AAG	CTC	CTG	CAC	GTC	GCC	AAG	GGC	CCG	CGC	GCC	CTG	CCC	1566
Phe	Arg	Gly	Lys	Leu	Leu	His	Val	Ala	Lys	Gly	Pro	Arg	Ala	Leu	Pro	
	295					300					305					
CAG	GGC	GCC	CCC	ACG	TCG	CCC	GGC	ATC	ACC	AAC	GCG	CTC	TGC	CT		

CTC GAC AAG CGG CTG TCC GCC CTC GCG AAG CGG CTG GGC TTC ACC TAC 1662  
 Leu Asp Lys Arg Leu Ser Ala Leu Ala Lys Arg Leu Gly Phe Thr Tyr  
 330 335 340

ACG CGC TAC GCG GAC GAC CTG ACC TTC TCC TGG ACG AAG GCG AAG CAG 1710  
 Thr Arg Tyr Ala Asp Asp Leu Thr Phe Ser Trp Thr Lys Ala Lys Gln  
 345 350 355

CCC AAG CCG CGG CGG ACG CAG CGT CCC CCC GTC GCG GTC CTC CTG TCT 1758  
 Pro Lys Pro Arg Arg Thr Gln Arg Pro Pro Val Ala Val Leu Leu Ser  
 360 365 370

CGC GTC CAG GAA GTG GTG GAG GCG GAG GGC TTC CGC GTG CAC CCG GAC 1806  
 Arg Val Gln Glu Val Val Glu Ala Glu Gly Phe Arg Val His Pro Asp  
 375 380 385

AAG ACG CGC GTC GCC CGC AAG GGC ACG CGG CAG CGG GTC ACC GGG CTC 1854  
 Lys Thr Arg Val Ala Arg Lys Gly Thr Arg Gln Arg Val Thr Gly Leu  
 390 395 400 405

GTC GTG AAT GCG GCG GGC AAG GAC GCG CCC GCG GCC CGA GTC CCG CGC 1902  
 Val Val Asn Ala Ala Gly Lys Asp Ala Pro Ala Ala Arg Val Pro Arg  
 410 415 420

GAC GTC GTC CGC CAG CTC CGC GGC GCC ATC CAC AAC CGG AAG AAG GGC 1950  
 Asp Val Val Arg Gln Leu Arg Ala Ala Ile His Asn Arg Lys Lys Gly  
 425 430 435

AAG CCG GGC CGC GAG GGC GAG TCG CTC GAG CAG CTC AAG GGC ATG GCC 1998  
 Lys Pro Gly Arg Glu Gly Glu Ser Leu Glu Gln Leu Lys Gly Met Ala  
 440 445 450

GCC TTC ATC CAC ATG ACG GAC CCG GCC AAG GGC CGC GCC TTC CTG GCT 2046  
 Ala Phe Ile His Met Thr Asp Pro Ala Lys Gly Arg Ala Phe Leu Ala  
 455 460 465

CAG CTC ACG GAG CTC GAG TCC ACG GCG AGC GCC GCT CCG CAG GCG GAG 2094  
 Gln Leu Thr Glu Leu Glu Ser Thr Ala Ser Ala Ala Pro Gln Ala Glu  
 470 475 480 485

TGACGCTCAG CGCGCGTCCG TCGCCGACGT GCCGCGCGCC AGCAACGCCG CATTGAGCAA 2154

CTCCGTCAGC CGGCGCGGGT AC 2176

(2) INFORMATION FOR SEQ ID NO:2:

(i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 263 amino acids
- (B) TYPE: amino acid
- (D) TOPOLOGY: linear

(ii) MOLECULE TYPE: protein

SECRET

[illegible]

(2) INFORMATION FOR SEQ ID NO:3:

(i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 263 amino acids
- (B) TYPE: amino acid
- (D) TOPOLOGY: linear

(ii) MOLECULE TYPE: protein

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:3:

Arg Pro Trp Ala Arg Thr Pro Pro Lys Ala Pro Arg Asn Gln Pro Val  
1 5 10 15  
Pro Phe Lys Pro Glu Arg Leu Gln Ala Leu Gln His Leu Val Arg Lys  
20 25 30  
Ala Leu Glu Ala Gly His Ile Glu Pro Tyr Thr Gly Pro Gly Asn Asn  
35 40 45  
Pro Val Phe Pro Val Lys Lys Ala Asn Gly Thr Trp Arg Phe Ile His  
50 55 60  
Asp Leu Arg Ala Thr Asn Ser Leu Thr Ile Asp Leu Ser Ser Ser Ser  
65 70 75 80  
Pro Gly Pro Pro Asp Leu Ser Ser Leu Pro Thr Thr Leu Ala His Leu  
85 90 95  
Gln Thr Ile Asp Leu Arg Asp Ala Phe Phe Gln Ile Pro Leu Pro Lys  
100 105 110  
Gln Phe Gln Pro Tyr Phe Ala Phe Thr Val Pro Gln Gln Cys Asn Tyr  
115 120 125  
Gly Pro Gly Thr Arg Tyr Ala Trp Lys Val Leu Pro Gln Gly Phe Lys  
130 135 140  
Asn Ser Pro Thr Leu Phe Glu Met Gln Leu Ala His Ile Leu Gln Pro  
145 150 155 160  
Ile Arg Gln Ala Phe Pro Gln Cys Thr Ile Leu Gln Tyr Met Asp Asp  
165 170 175  
Ile Leu Leu Ala Ser Pro Ser His Glu Asp Leu Leu Leu Leu Ser Glu  
180 185 190  
Ala Thr Met Ala Ser Leu Ile Ser His Gly Leu Pro Val Ser Glu Asn  
195 200 205  
Lys Thr Gln Gln Thr Pro Gly Thr Ile Lys Phe Leu Gly Gln Ile Ile  
210 215 220

SECRET  
Sub  
Gmt



Trp Val Ser Lys Gly Thr Pro  
260

Tyr Gly Cys Thr Tyr Ser Arg Tyr Ala Asp Asp Ile Thr Ile Ser Thr  
165 170 175

RECEIVED  
JUN 10 1964  
Sub  
G  
cont

Asn Lys Asn Thr Phe Pro Leu Glu Met Ala Thr Val Gln Pro Glu Gly  
 180 185 190  
 Val Val Leu Gly Lys Val Leu Val Lys Glu Ile Glu Asn Ser Gly Phe  
 195 200 205  
 Glu Ile Asn Asp Ser Lys Thr Arg Leu Thr Tyr Lys Thr Ser Arg Gln  
 210 215 220  
 Glu Val Thr Gly Leu Thr Val Asn Arg Ile Val Asn Ile Asp Arg Cys  
 225 230 235 240  
 Tyr Tyr Lys Lys Thr Arg Ala Leu Ala His Ala Leu Tyr Arg Thr Gly  
 245 250 255  
 Glu Tyr Lys

(2) INFORMATION FOR SEQ ID NO:5:

(i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 266 amino acids
- (B) TYPE: amino acid
- (D) TOPOLOGY: linear

(ii) MOLECULE TYPE: protein

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:5:

Ala Phe His Arg Glu Val Asp Thr Ala Thr His Tyr Val Ser Trp Thr  
 1 5 10 15  
 Ile Pro Lys Arg Asp Gly Ser Lys Arg Thr Ile Thr Ser Pro Lys Pro  
 20 25 30  
 Glu Leu Lys Ala Ala Gln Arg Trp Val Leu Ser Asn Val Val Glu Arg  
 35 40 45  
 Leu Pro Val His Gly Ala Ala His Gly Phe Val Ala Gly Arg Ser Ile  
 50 55 60  
 Leu Thr Asn Ala Leu Ala His Gln Gly Ala Asp Val Val Val Lys Val  
 65 70 75 80  
 Asp Leu Lys Asp Phe Phe Pro Ser Val Thr Trp Arg Arg Val Lys Gly  
 85 90 95  
 Leu Leu Arg Lys Gly Gly Leu Arg Glu Gly Thr Ser Thr Leu Leu Ser  
 100 105 110  
 Leu Leu Ser Thr Glu Ala Pro Arg Glu Ala Val Gln Phe Arg Gly Lys  
 115 120 125

08808001-00000000

cont

Leu Leu His Val Ala Lys Gly Pro Arg Ala Leu Pro Gln Gly Ala Pro  
 130 135 140  
 Thr Ser Pro Gly Ile Thr Asn Ala Leu Cys Leu Lys Leu Asp Lys Arg  
 145 150 155 160  
 Leu Ser Ala Leu Ala Lys Arg Leu Gly Phe Thr Tyr Thr Arg Tyr Ala  
 165 170 175  
 Asp Asp Leu Thr Phe Ser Trp Thr Lys Ala Lys Gln Pro Lys Pro Arg  
 180 185 190  
 Arg Thr Gln Arg Pro Pro Val Ala Val Leu Leu Ser Arg Val Gln Glu  
 195 200 205  
 Val Val Glu Ala Glu Gly Phe Arg Val His Pro Asp Lys Thr Arg Val  
 210 215 220  
 Ala Arg Lys Gly Thr Arg Gln Arg Val Thr Gly Leu Val Val Asn Ala  
 225 230 235 240  
 Ala Gly Lys Asp Ala Pro Ala Ala Arg Val Pro Arg Asp Val Val Arg  
 245 250 255  
 Gln Leu Arg Ala Ala Ile His Asn Arg Lys  
 260 265

(2) INFORMATION FOR SEQ ID NO:6:

(i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 111 amino acids
- (B) TYPE: amino acid
- (D) TOPOLOGY: linear

(ii) MOLECULE TYPE: protein

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:6:

Pro Thr Pro Glu Leu Thr Ala Pro Ser Ser Asp Ala Ala Ala Lys Arg  
 1 5 10 15  
 Glu Ala Arg Arg Leu Ala His Glu Ala Leu Leu Val Arg Ala Lys Ala  
 20 25 30  
 Ile Asp Glu Ala Gly Gly Ala Asp Asp Trp Val Gln Ala Gln Leu Val  
 35 40 45  
 Ser Lys Gly Leu Ala Val Glu Asp Leu Asp Phe Ser Ser Ala Ser Glu  
 50 55 60  
 Lys Asp Lys Lys Ala Trp Lys Glu Lys Lys Lys Ala Glu Ala Thr Glu  
 65 70 75 80

0308001-0308001

Seq  
G  
cont

Arg Arg Ala Leu Lys Arg Gln Ala His Glu Ala Trp Lys Ala Thr His  
85 90 95

Val Gly His Leu Gly Ala Gly Val His Trp Ala Glu Asp Arg Leu  
100 105 110

(2) INFORMATION FOR SEQ ID NO:7:

(i) SEQUENCE CHARACTERISTICS:

(A) LENGTH: 110 amino acids

(B) TYPE: amino acid

(D) TOPOLOGY: linear

(ii) MOLECULE TYPE: protein

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:7:

Pro Asp Pro Asp Met Thr Arg Val Thr Asn Ser Pro Ser Leu Gln Ala  
1 5 10 15

His Leu Gln Ala Leu Tyr Leu Val Gln His Glu Val Trp Arg Pro Leu  
20 25 30

Ala Ala Ala Tyr Gln Glu Gln Leu Asp Arg Pro Val Val Pro His Pro  
35 40 45

Tyr Arg Val Gly Asp Thr Val Trp Val Arg Arg His Gln Thr Lys Asn  
50 55 60

Leu Glu Pro Arg Trp Lys Gly Pro Tyr Thr Val Leu Leu Thr Thr Pro  
65 70 75 80

Thr Ala Leu Lys Val Asp Gly Ile Ala Ala Trp Ile His Ala Ala His  
85 90 95

Val Lys Ala Ala Asp Pro Gly Gly Gly Pro Ser Ser Arg Leu  
100 105 110

(2) INFORMATION FOR SEQ ID NO:8:

(i) SEQUENCE CHARACTERISTICS:

(A) LENGTH: 75 amino acids

(B) TYPE: amino acid

(D) TOPOLOGY: linear

(ii) MOLECULE TYPE: protein

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:8:

Gly Lys Asp Ala Pro Ala Ala Arg Val Pro Arg Asp Val Val Arg Gln

RECEIVED  
JUN 10 1980  
FBI  
cont





(xi) SEQUENCE DESCRIPTION: SEQ ID NO:12:

Tyr Lys Asn Leu Leu Pro Gln Gly Ala Pro Ser Ser Pro Lys Leu Ala  
1 5 10 15  
Asn Leu Ile Cys Ser Lys Leu Asp Tyr Arg Ile Gln Gly Tyr Ala Gly  
20 25 30  
Ser Arg Gly Leu Ile Tyr Thr Arg Tyr Ala Asp Asp Leu Thr Leu Ser  
35 40 45  
Ala Gln Ser Met Lys Lys Val Val Lys Ala Arg Asp Phe Leu Phe Ser  
50 55 60  
Ile Ile Pro Ser  
65

(2) INFORMATION FOR SEQ ID NO:13:

(i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 67 amino acids  
(B) TYPE: amino acid  
(D) TOPOLOGY: linear

(ii) MOLECULE TYPE: protein

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:13:

Tyr Gln Tyr Asn Val Leu Pro Gln Gly Trp Lys Gly Ser Pro Ala Ile  
1 5 10 15  
Phe Gln Ser Ser Met Thr Lys Ile Leu Glu Pro Phe Lys Lys Gln Asn  
20 25 30  
Pro Asp Ile Val Ile Tyr Gln Tyr Met Asp Asp Leu Tyr Val Gly Ser  
35 40 45  
Asp Leu Glu Ile Gly Gln His Arg Thr Lys Ile Glu Glu Leu Arg Gln  
50 55 60  
His Leu Leu  
65

(2) INFORMATION FOR SEQ ID NO:14:

(i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 66 amino acids  
(B) TYPE: amino acid  
(D) TOPOLOGY: linear

(ii) MOLECULE TYPE: protein

RECEIVED 1-03-03  
Sub  
G-1  
cont

**SECRET**

(2) INFORMATION FOR SEQ ID NO:15:

(A) LENGTH: 65 amino acids  
(B) TYPE: amino acid  
(D) TOPOLOGY: linear

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:15:

(2) INFORMATION FOR SEQ ID NO:16:

(A) LENGTH: 65 amino acids  
(B) TYPE: amino acid  
(D) TOPOLOGY: linear

(ii) MOLECULE TYPE: protein



[illegible]

100

100

100

- 100

100

100

100

100

100

- 100



00000000000000000000

Tyr 1	Glu	Phe	Cys 5	Arg	Leu	Pro	Phe	Gly	Leu 10	Arg	Asn	Ala	Ser	Ser 15	Ile
Phe	Gln	Arg	Ala 20	Leu	Asp	Asp	Val	Leu 25	Arg	Glu	Gln	Ile	Gly 30	Lys	Ile
Cys	Tyr	Val 35	Tyr	Val	Asp	Asp	Val 40	Ile	Ile	Phe	Ser	Glu 45	Asn	Glu	Ser
Asp	His 50	Val	Arg	His	Ile	Asp 55	Thr	Val	Leu	Lys	Cys 60	Leu			

(2) INFORMATION FOR SEQ ID NO:21:

(i) SEQUENCE CHARACTERISTICS:  
(A) LENGTH: 63 amino acids  
(B) TYPE: amino acid  
(D) TOPOLOGY: linear

(ii) MOLECULE TYPE: protein

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:21:

Cys	Lys	Leu	Asn	Lys	Ala	Ile	Tyr	Gly	Leu	Lys	Gln	Ala	Ala	Arg	Cys
1				5					10					15	
Trp	Phe	Arg	Cys	Ile	Tyr	Ile	Leu	Asp	Lys	Gly	Asn	Ile	Asn	Glu	Asn
			20					25					30		
Ile	Tyr	Val	Leu	Leu	Tyr	Val	Asp	Asp	Val	Val	Ile	Ala	Thr	Gly	Asp
		35					40					45			
Met	Thr	Arg	Met	Asn	Asn	Phe	Lys	Arg	Tyr	Leu	Met	Glu	Lys	Phe	
	50					55					60				

(2) INFORMATION FOR SEQ ID NO:22:

(i) SEQUENCE CHARACTERISTICS:  
 (A) LENGTH: 62 amino acids  
 (B) TYPE: amino acid  
 (D) TOPOLOGY: linear

(ii) MOLECULE TYPE: protein

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:22:

Cys Leu Leu Lys Lys Ser Leu Tyr Gly Leu Lys Gln Ser Pro Arg Gln  
1 5 10 15



(D) OTHER INFORMATION: /note= "This region can hydrogen bond to nucleotides 61-67 of SEQ ID NO: 25 of this application."

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:24:

CACGCAUGUA GGCAGAUUUG UUGGUUGUGA AUCGCAACCA GUGGCCUUAU UGGCAGGA

58

(2) INFORMATION FOR SEQ ID NO:25:

(i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 67 base pairs
- (B) TYPE: nucleic acid
- (C) STRANDEDNESS: single
- (D) TOPOLOGY: linear

(ix) FEATURE:

- (A) NAME/KEY: misc\_feature
- (B) LOCATION: 1
- (D) OTHER INFORMATION: /note= "The 5' position of this nucleotide is linked to the 2' position of nucleotide number 15 of SEQ ID NO: 24 of this application."

(ix) FEATURE:

- (A) NAME/KEY: misc\_binding
- (B) LOCATION: 61..67
- (D) OTHER INFORMATION: /note= "This region can hydrogen bond to nucleotides 52-58 of SEQ ID NO: 24 of this application."

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:25:

TCCTTCGCAC AGCACACCTG CCGTATAGCT CTGAATCAAG GATTTTAGGG AGGCGATTCC

60

TCCTGCC

67

(2) INFORMATION FOR SEQ ID NO:26:

(i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 2423 base pairs
- (B) TYPE: nucleic acid
- (C) STRANDEDNESS: double
- (D) TOPOLOGY: linear

(ix) FEATURE:

- (A) NAME/KEY: CDS
- (B) LOCATION: 418..2175

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:26:

TGGCCATTNA GATACGGATT TTCACTTCCT TGACAGTGCA TGA CTATGCT GCATGAAATN 60  
GCATGATCGA TTGAGGATCG TCTTTGCTCA GATCCGCCAG AACTGGCGGG CTTTGTCTCA 120  
TGTCATGCAT GTGCATGAAA ACCACTGCAT AAAGCGGGCA GCGTGGCGG GGATACGAGC 180  
GCGCGCTATC ACCGAAAATA GCCAAAATAC TTCTGGAAAA CAGAAAGTTG AAGTGATATG 240  
TTCATAAACA CGCATGTAGG CAGATTTGTT GGTGTGAAT CGCAACCAGT GGCCTTAATG 300  
GCAGGAGGAA TCGCCTCCCT AAAATCCTTG ATTCAGAGCT ATACGGCAGG TGTGCTGTGC 360  
GAAGGAGTGC CTGCATGCGT TTCTCCTTGG CCTTTTTTCC TCTGGGATGA AGAAGAA 417  
ATG ACA AAA ACA TCT AAA CTT GAC GCA CTT AGG GCT GCT ACT TCA CGT 465  
Met Thr Lys Thr Ser Lys Leu Asp Ala Leu Arg Ala Ala Thr Ser Arg  
1 5 10 15  
GAA GAC TTG GCT AAA ATT TTA GAT ATT AAG TTG GTA TTT TTA ACT AAC 513  
Glu Asp Leu Ala Lys Ile Leu Asp Ile Lys Leu Val Phe Leu Thr Asn  
20 25 30  
GTT CTA TAT AGA ATC GGC TCG GAT AAT CAA TAC ACT CAA TTT ACA ATA 561  
Val Leu Tyr Arg Ile Gly Ser Asp Asn Gln Tyr Thr Gln Phe Thr Ile  
35 40 45  
CCG AAG AAA GGA AAA GGG GTA AGG ACT ATT TCT GCA CCT ACA GAC CGG 609  
Pro Lys Lys Gly Lys Gly Val Arg Thr Ile Ser Ala Pro Thr Asp Arg  
50 55 60  
TTG AAG GAC ATC CAA CGA AGA ATA TGT GAC TTA CTT TCT GAT TGT AGA 657  
Leu Lys Asp Ile Gln Arg Arg Ile Cys Asp Leu Leu Ser Asp Cys Arg  
65 70 75 80  
GAT GAG ATC TTT GCT ATA AGG AAA ATT AGT AAC AAC TAT TCC TTT GGT 705  
Asp Glu Ile Phe Ala Ile Arg Lys Ile Ser Asn Asn Tyr Ser Phe Gly  
85 90 95  
TTT GAG AGG GGA AAA TCA ATA ATC CTA AAT GCT TAT AAG CAT AGA GGC 753  
Phe Glu Arg Gly Lys Ser Ile Ile Leu Asn Ala Tyr Lys His Arg Gly  
100 105 110  
AAA CAA ATA ATA TTA AAT ATA GAT CTT AAG GAT TTT TTT GAA AGC TTT 801  
Lys Gln Ile Ile Leu Asn Ile Asp Leu Lys Asp Phe Phe Glu Ser Phe  
115 120 125  
AAT TTT GGA CGA GTT AGA GGA TAT TTT CTT TCC AAT CAG GAT TTT TTA 849  
Asn Phe Gly Arg Val Arg Gly Tyr Phe Leu Ser Asn Gln Asp Phe Leu  
130 135 140  
TTA AAT CCT GTG GTG GCA ACG ACA CTT GCA AAA GCT GCA TGC TAT AAT 897  
Leu Asn Pro Val Val Ala Thr Thr Leu Ala Lys Ala Ala Cys Tyr Asn

145				150				155				160				
GGA Gly	ACC Thr	CTC Leu	CCC Pro	CAA Gln 165	GGA Gly	AGT Ser	CCA Pro	TGT Cys	TCT Ser 170	CCT Pro	ATT Ile	ATC Ile	TCA Ser	AAT Asn 175	CTA Leu	945
ATT Ile	TGC Cys	AAT Asn	ATT Ile 180	ATG Met	GAT Asp	ATG Met	AGA Arg	TTA Leu 185	GCT Ala	AAG Lys	CTG Leu	GCT Ala	AAA Lys 190	AAA Lys	TAT Tyr	993
GGA Gly	TGT Cys	ACT Thr 195	TAT Tyr	AGC Ser	AGA Arg	TAT Tyr	GCT Ala 200	GAT Asp	GAT Asp	ATA Ile	ACA Thr	ATT Ile 205	TCT Ser	ACA Thr	AAT Asn	1041
AAA Lys	AAT Asn 210	ACA Thr	TTT Phe	CCG Pro	TTA Leu	GAA Glu 215	ATG Met	GCT Ala	ACT Thr	GTG Val	CAA Gln 220	CCT Pro	GAA Glu	GGG Gly	GTT Val	1089
GTT Val 225	TTG Leu	GGA Gly	AAA Lys	GTT Val	TTG Leu 230	GTA Val	AAA Lys	GAA Glu	ATA Ile	GAA Glu 235	AAC Asn	TCT Ser	GGA Gly	TTC Phe	GAA Glu 240	1137
ATA Ile	AAT Asn	GAT Asp	TCA Ser	AAG Lys 245	ACT Thr	AGG Arg	CTT Leu	ACG Thr	TAT Tyr 250	AAG Lys	ACA Thr	TCA Ser	AGG Arg	CAA Gln 255	GAA Glu	1185
GTA Val	ACG Thr	GGA Gly	CTT Leu 260	ACA Thr	GTT Val	AAC Asn	AGA Arg	ATC Ile 265	GTT Val	AAT Asn	ATT Ile	GAT Asp	AGA Arg 270	TGT Cys	TAT Tyr	1233
TAT Tyr	AAA Lys	AAA Lys	ACT Thr 275	CGG Arg	GCG Ala	TTG Leu	GCA Ala 280	CAT His	GCT Ala	TTG Leu	TAT Tyr	CGT Arg 285	ACA Thr	GGT Gly	GAA Glu	1281
TAT Tyr	AAA Lys 290	GTG Val	CCA Pro	GAT Asp	GAA Glu	AAT Asn 295	GGT Gly	GTT Val	TTA Leu	GTT Val	TCA Ser 300	GGA Gly	GGT Gly	CTG Leu	GAT Asp	1329
AAA Lys 305	CTT Leu	GAG Glu	GGG Gly	ATG Met 310	TTT Phe	GGT Gly	TTT Phe	ATT Ile	GAT Asp	CAA Gln 315	GTT Val	GAT Asp	AAG Lys	TTT Phe 320	AAC Asn	1377
AAT Asn	ATA Ile	AAG Lys	AAA Lys 325	AAA Lys	CTG Leu	AAC Asn	AAG Lys	CAA Gln	CCT Pro 330	GAT Asp	AGA Arg	TAT Tyr	GTA Val	TTG Leu 335	ACT Thr	1425
AAT Asn	GCG Ala	ACT Thr 340	TTG Leu	CAT His	GGT Gly	TTT Phe	AAA Lys	TTA Leu 345	AAG Lys	TTG Leu	AAT Asn	GCG Ala	CGA Arg 350	GAA Glu	AAA Lys	1473
GCA Ala	TAT Tyr	AGT Ser 355	AAA Lys	TTT Phe	ATT Ile	TAC Tyr	TAT Tyr 360	AAA Lys	TTT Phe	TTT Phe	CAT His	GGC Gly 365	AAC Asn	ACC Thr	TGT Cys	1521

CCT Pro	ACG Thr	ATA Ile	ATT Ile	ACA Thr	GAA Glu	GGG Gly	AAG Lys	ACT Thr	GAT Asp	CGG Arg	ATA Ile	TAT Tyr	TTG Leu	AAG Lys	GCT Ala	1569
370 375 380																
GCT Ala	TTG Leu	CAT His	TCT Ser	TTG Leu	GAG Glu	ACA Thr	TCA Ser	TAT Tyr	CCT Pro	GAG Glu	TTG Leu	TTT Phe	AGA Arg	GAA Glu	AAA Lys	1617
385 390 395 400																
ACA Thr	GAT Asp	AGT Ser	AAA Lys	AAG Lys	AAA Lys	GAA Glu	ATA Ile	AAT Asn	CTT Leu	AAT Asn	ATA Ile	TTT Phe	AAA Lys	TCT Ser	AAT Asn	1665
405 410 415																
GAA Glu	AAG Lys	ACC Thr	AAA Lys	TAT Tyr	TTT Phe	TTA Leu	GAT Asp	CTT Leu	TCT Ser	GGG Gly	GGA Gly	ACT Thr	GCA Ala	GAT Asp	CTG Leu	1713
420 425 430																
AAA Lys	AAA Lys	TTT Phe	GTA Val	GAG Glu	CGT Arg	TAT Tyr	AAA Lys	AAT Asn	AAT Asn	TAT Tyr	GCT Ala	TCT Ser	TAT Tyr	TAT Tyr	GGT Gly	1761
435 440 445																
TCT Ser	GTT Val	CCA Pro	AAA Lys	CAG Gln	CCA Pro	GTG Val	ATT Ile	ATG Met	GTT Val	CTT Leu	GAT Asp	AAT Asn	GAT Asp	ACA Thr	GGT Gly	1809
450 455 460																
CCA Pro	AGC Ser	GAT Asp	TTA Leu	CTT Leu	AAT Asn	TTT Phe	CTG Leu	CGC Arg	AAT Asn	AAA Lys	GTT Val	AAA Lys	AGC Ser	TGC Cys	CCA Pro	1857
465 470 475 480																
GAC Asp	GAT Asp	GTA Val	ACT Thr	GAA Glu	ATG Met	AGA Arg	AAG Lys	ATG Met	AAA Lys	TAT Tyr	ATT Ile	CAT His	GTT Val	TTC Phe	TAT Tyr	1905
485 490 495																
AAT Asn	TTA Leu	TAT Tyr	ATA Ile	GTT Val	CTC Leu	ACA Thr	CCA Pro	TTG Leu	AGT Ser	CCT Pro	TCC Ser	GGC Gly	GAA Glu	CAA Gln	ACT Thr	1953
500 505 510																
TCA Ser	ATG Met	GAG Glu	GAT Asp	CTT Leu	TTC Phe	CCT Pro	AAA Lys	GAT Asp	ATT Ile	TTA Leu	GAT Asp	ATC Ile	AAG Lys	ATT Ile	GAT Asp	2001
515 520 525																
GGT Gly	AAG Lys	AAA Lys	TTC Phe	AAC Asn	AAA Lys	AAT Asn	AAT Asn	GAT Asp	GGA Gly	GAC Asp	TCA Ser	AAA Lys	ACG Thr	GAA Glu	TAT Tyr	2049
530 535 540																
GGG Gly	AAG Lys	CAT His	ATT Ile	TTT Phe	TCC Ser	ATG Met	AGG Arg	GTT Val	GTT Val	AGA Arg	GAT Asp	AAA Lys	AAG Lys	CGG Arg	AAA Lys	2097
545 550 555 560																
ATA Ile	GAT Asp	TTT Phe	AAG Lys	GCA Ala	TTT Phe	TGT Cys	TGT Cys	ATT Ile	TTT Phe	GAT Asp	GCT Ala	ATA Ile	AAA Lys	GAT Asp	ATA Ile	2145
565 570 575																
AAG Lys	GAA Glu	CAT His	TAT Tyr	AAA Lys	TTA Leu	ATG Met	TTA Leu	AAT Asn	AGC Ser	TAATGAACAG	CCCTAACGTT					2195







Val Thr Asn Lys Gly Arg Gln Lys Val Val Pro Leu Thr Asn Thr Thr  
450 455 460

Asn Gln Lys Thr Glu Leu Gln Ala Ile Tyr Leu Ala Leu Gln Asp Ser  
465 470 475 480

Gly Leu Glu Val Asn Ile Val Thr Asp Ser Gln Tyr Ala Leu Gln Ile  
485 490 495

Ile Gln Ala Gln Pro Asp Lys Ser Glu Ser Glu Leu Val Asn Gln Ile  
500 505 510

Ile Glu Gln Leu Ile Lys Lys Glu Lys Val Tyr Leu Ala Trp Val Pro  
515 520 525

Ala His Lys Gly Ile Gly Gly Asn Glu Gln Val Asp Lys Leu Val Ser  
530 535 540

Ala Gly  
545

(2) INFORMATION FOR SEQ ID NO:28:

(i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 578 amino acids
- (B) TYPE: amino acid
- (D) TOPOLOGY: linear

(ii) MOLECULE TYPE: protein

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:28:

Arg Pro Trp Ala Arg Thr Pro Pro Lys Ala Pro Arg Asn Gln Pro Val  
1 5 10 15

Pro Phe Lys Pro Glu Arg Leu Gln Ala Leu Gln His Leu Val Arg Lys  
20 25 30

Ala Leu Glu Ala Gly His Ile Glu Pro Tyr Thr Gly Pro Gly Asn Asn  
35 40 45

Pro Val Phe Pro Val Lys Lys Ala Asn Gly Thr Trp Arg Phe Ile His  
50 55 60

Asp Leu Arg Ala Thr Asn Ser Leu Thr Ile Asp Leu Ser Ser Ser Ser  
65 70 75 80

Pro Gly Pro Pro Asp Leu Ser Ser Leu Pro Thr Thr Leu Ala His Leu  
85 90 95

Gln Thr Ile Asp Leu Arg Asp Ala Phe Phe Gln Ile Pro Leu Pro Lys  
100 105 110

1-008001-030397







4500 T. 030000

(i) SEQUENCE CHARACTERISTICS:

- (ii) MOLECULE TYPE: protein

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:30:

Arg Trp Phe Ser Phe His Arg Glu Val Asp Thr Gly Thr His Tyr Gln  
1 5 10 15  
Thr Trp Glu Ile Pro Lys Arg Asp Gly Gly Lys Arg Thr Leu Thr Ala  
20 25 30  
Pro Lys Arg Glu Leu Lys Ala Val Gln Arg Trp Val Leu Ala Asn Val  
35 40 45  
Val Glu Arg Leu Pro Val His Gly Ala Ala His Gly Phe Val Ala Gly  
50 55 60  
Arg Ser Ile Leu Thr Asn Ala Leu Ala His Gln Gly Ala Asp Val Val  
65 70 75 80  
Val Lys Val Asp Met Lys Asp Phe Phe Pro Ser Val Thr Trp Pro Arg  
85 90 95  
Val Lys Gly Leu Leu Arg Lys Gly Gly Leu Pro Glu Asn Leu Ala Thr  
100 105 110  
Leu Leu Ala Leu Leu Ser Thr Glu Ala Pro Arg Glu Val Val Arg Phe  
115 120 125  
Arg Gly Glu Thr Leu Tyr Val Ala Lys Gly Pro Arg Ala Leu Pro Gln  
130 135 140  
Gly Ala Pro Thr Ser Pro Ala Leu Thr Asn Ala Leu Cys Leu Arg Leu  
145 150 155 160  
Asp Lys Arg Leu Ser Ala Leu Ser Lys Arg Leu Gly Phe Thr Tyr Thr  
165 170 175  
Arg Tyr Ala Asp Asp Leu Thr Phe Ser Trp Arg Arg Ala Lys Lys Ser  
180 185 190  
Arg Gln Lys Glu Leu Pro Leu Ala Asp Ala Pro Val Ala Leu Leu Leu  
195 200 205  
Ala Arg Val Lys Gly Val Leu Glu Ala Glu Gly Phe Thr Leu His Pro  
210 215 220  
Asp Lys Thr Arg Val Gln Arg Lys Gly Ser Arg Gln Arg Val Thr Gly  
225 230 235 240  
Leu Val Val

(2) INFORMATION FOR SEQ ID NO:31:

(i) SEQUENCE CHARACTERISTICS:



(A) LENGTH: 241 amino acids  
(B) TYPE: amino acid  
(D) TOPOLOGY: linear

(ii) MOLECULE TYPE: protein

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:31:

Arg Trp Phe Ala Phe His Arg Glu Val Asp Thr Ala Thr His Tyr Val  
1 5 10 15  
Ser Trp Thr Ile Pro Lys Arg Asp Gly Ser Lys Arg Thr Ile Thr Ser  
20 25 30  
Pro Lys Pro Glu Leu Lys Ala Ala Gln Arg Trp Val Leu Ser Asn Val  
35 40 45  
Val Glu Arg Leu Pro Val His Gly Ala Ala His Gly Phe Val Ala Gly  
50 55 60  
Arg Ser Ile Leu Thr Asn Ala Leu Ala His Gln Gly Ala Asp Val Val  
65 70 75 80  
Val Lys Val Asp Leu Lys Asp Phe Phe Pro Ser Val Thr Trp Arg Arg  
85 90 95  
Val Lys Gly Leu Leu Arg Lys Gly Gly Leu Arg Glu Gly Thr Ser Thr  
100 105 110  
Leu Leu Ser Leu Leu Ser Thr Glu Ala Pro Arg Glu Ala Val Gln Phe  
115 120 125  
Pro Arg Glu Leu Leu His Val Ala Lys Gly Pro Arg Ala Leu Pro Gln  
130 135 140  
Gly Ala Pro Thr Ser Pro Gly Ile Thr Asn Ala Leu Cys Leu Lys Leu  
145 150 155 160  
Asp Lys Arg Leu Ser Ala Leu Ala Lys Arg Leu Gly Phe Thr Tyr Thr  
165 170 175  
Arg Tyr Ala Asp Asp Leu Thr Phe Ser Trp Thr Lys Ala Lys Gln Pro  
180 185 190  
Lys Pro Arg Arg Thr Gln Arg Pro Pro Val Ala Val Leu Leu Ser Arg  
195 200 205  
Val Gln Glu Val Val Glu Ala Glu Gly Phe Arg Val His Pro Asp Lys  
210 215 220  
Thr Arg Val Ala Arg Lys Gly Thr Arg Gln Arg Val Thr Gly Leu Val  
225 230 235 240

RECEIVED FEB 27 1997

Val

(2) INFORMATION FOR SEQ ID NO:32:

(i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 231 amino acids
- (B) TYPE: amino acid
- (D) TOPOLOGY: linear

(ii) MOLECULE TYPE: protein

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:32:

Arg His Tyr Ser Ile His Arg Pro Arg Glu Arg Val Arg His Tyr Val  
1 5 10 15  
Thr Phe Ala Val Pro Lys Arg Ser Gly Gly Val Arg Leu Leu His Ala  
20 25 30  
Pro Lys Arg Arg Leu Lys Ala Leu Gln Arg Arg Met Leu Ala Leu Leu  
35 40 45  
Val Ser Lys Leu Pro Val Ser Pro Gln Ala His Gly Phe Val Pro Gly  
50 55 60  
Arg Ser Ile Lys Thr Gly Ala Ala Pro His Val Gly Arg Arg Val Val  
65 70 75 80  
Leu Lys Leu Asp Leu Lys Asp Phe Phe Pro Ser Val Thr Phe Ala Arg  
85 90 95  
Val Arg Gly Leu Leu Lys Ala Leu Gly Tyr Gly Tyr Pro Val Ala Ala  
100 105 110  
Thr Leu Ala Val Leu Met Thr Glu Ser Glu Arg Gln Pro Val Glu Leu  
115 120 125  
Glu Gly Ile Leu Phe His Val Pro Val Gly Pro Arg Val Cys Val Gln  
130 135 140  
Gly Ala Pro Thr Ser Pro Ala Leu Cys Asn Ala Val Leu Leu Arg Leu  
145 150 155 160  
Asp Arg Arg Leu Ala Gly Leu Ala Arg Arg Tyr Gly Tyr Thr Tyr Thr  
165 170 175  
Arg Tyr Ala Asp Asp Leu Thr Phe Ser Gly Asp Asp Val Thr Ala Leu  
180 185 190  
Glu Arg Val Arg Ala Leu Ala Ala Arg Tyr Val Gln Glu Glu Gly Phe  
195 200 205

0808031.030397

Glu Val Asn Arg Glu Lys Thr Arg Val Gln Arg Arg Gly Gly Ala Gln  
 210 215 220

Arg Val Thr Gly Val Thr Val  
 225 230

(2) INFORMATION FOR SEQ ID NO:33:

(i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 234 amino acids
- (B) TYPE: amino acid
- (D) TOPOLOGY: linear

(ii) MOLECULE TYPE: protein

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:33:

Phe Leu Thr Asn Val Leu Tyr Arg Ile Gly Ser Asp Asn Gln Tyr Thr  
 1 5 10 15  
 Gln Phe Thr Ile Pro Lys Lys Gly Lys Gly Val Arg Thr Ile Ser Ala  
 20 25 30  
 Pro Thr Asp Arg Leu Lys Asp Ile Gln Arg Arg Ile Cys Asp Leu Leu  
 35 40 45  
 Ser Asp Cys Arg Asp Glu Ile Phe Ala Ile Arg Lys Ile Ser Asn Asn  
 50 55 60  
 Tyr Ser Phe Gly Phe Glu Arg Gly Lys Ser Ile Ile Leu Asn Ala Tyr  
 65 70 75 80  
 Lys His Arg Gly Lys Gln Ile Ile Leu Asn Ile Asp Leu Lys Asp Phe  
 85 90 95  
 Phe Glu Ser Phe Asn Phe Gly Arg Val Arg Gly Tyr Phe Leu Ser Asn  
 100 105 110  
 Gln Asp Phe Leu Leu Asn Pro Val Val Ala Thr Thr Leu Ala Lys Ala  
 115 120 125  
 Ala Cys Tyr Asn Gly Thr Leu Pro Gln Gly Ser Pro Cys Ser Pro Ile  
 130 135 140  
 Ile Ser Asn Leu Ile Cys Asn Ile Met Asp Met Arg Leu Ala Lys Leu  
 145 150 155 160  
 Ala Lys Lys Tyr Gly Cys Thr Tyr Ser Arg Tyr Ala Asp Asp Ile Thr  
 165 170 175  
 Ile Ser Thr Asn Lys Asn Thr Phe Pro Leu Glu Met Ala Thr Val Gln  
 180 185 190

0808031-030397

Pro Glu Gly Val Val Leu Gly Lys Val Leu Val Lys Glu Ile Glu Asn  
195 200 205

Ser Gly Phe Glu Ile Asn Asp Ser Lys Thr Arg Leu Thr Tyr Lys Thr  
210 215 220

Ser Arg Gln Glu Val Thr Gly Leu Thr Val  
225 230

(2) INFORMATION FOR SEQ ID NO:34:

(i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 215 amino acids
- (B) TYPE: amino acid
- (D) TOPOLOGY: linear

(ii) MOLECULE TYPE: protein

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:34:

Val Glu Thr Leu Arg Leu Leu Ile Tyr Thr Ala Asp Phe Arg Tyr Arg  
1 5 10 15

Ile Tyr Thr Val Glu Lys Lys Gly Pro Glu Lys Arg Met Arg Thr Ile  
20 25 30

Tyr Gln Pro Ser Arg Glu Leu Lys Ala Leu Gln Gly Trp Val Leu Arg  
35 40 45

Asn Ile Leu Asp Lys Leu Ser Ser Ser Pro Phe Ser Ile Gly Phe Glu  
50 55 60

Lys His Gln Ser Ile Leu Asn Asn Ala Thr Pro His Ile Gly Ala Asn  
65 70 75 80

Phe Ile Leu Asn Ile Asp Leu Glu Asp Phe Phe Pro Ser Leu Thr Ala  
85 90 95

Asn Lys Val Phe Gly Val Phe His Ser Leu Gly Tyr Asn Arg Leu Ile  
100 105 110

Ser Ser Val Leu Thr Lys Ile Cys Cys Tyr Lys Asn Leu Leu Pro Gln  
115 120 125

Gly Ala Pro Ser Ser Pro Lys Leu Ala Asn Leu Ile Cys Ser Lys Leu  
130 135 140

Asp Tyr Arg Ile Gln Gly Tyr Ala Gly Ser Arg Gly Leu Ile Tyr Thr  
145 150 155 160

Arg Tyr Ala Asp Asp Leu Thr Leu Ser Ala Gln Ser Met Lys Lys Val  
165 170 175

DEB0341-030397



069070

(i) SEQUENCE CHARACTERISTICS:

- (ii) MOLECULE TYPE: protein

Ile 1	Gln	Arg	Leu 5	His	Ala	Leu	Ser	Asn 10	His	Ala	Gly	Arg	His	Tyr 15	Arg
Arg	Ile	Ile 20	Leu	Ser	Lys	Arg	His	Gly 25	Gly	Gln	Arg	Leu 30	Val	Leu	Ala
Pro	Asp 35	Tyr	Leu	Leu	Lys	Thr	Val 40	Gln	Arg	Asn	Ile 45	Leu	Lys	Asn	Val
Leu	Ser 50	Gln	Phe	Pro	Leu	Ser 55	Pro	Phe	Ala	Thr	Ala 60	Tyr	Arg	Pro	Gly
Cys 65	Pro	Ile	Val	Ser 70	Asn	Ala	Gln	Pro	His	Cys 75	Gln	Gln	Pro	Gln	Ile 80
Leu	Lys	Leu	Asp 85	Ile	Glu	Asn	Phe	Phe	Asp 90	Ser	Ile	Ser	Trp	Leu 95	Gln
Val	Trp	Arg	Val 100	Phe	Arg	Gln	Ala	Gln	Leu 105	Pro	Arg	Asn	Val	Val	Thr
Met	Leu 115	Thr	Trp	Ile	Cys	Cys	Tyr 120	Asn	Asp	Ala	Leu	Pro 125	Gln	Gly	Ala
Pro	Thr 130	Ser	Pro	Ala	Ile	Ser 135	Asn	Leu	Val	Met	Arg 140	Arg	Phe	Asp	Glu
Arg 145	Ile	Gly	Glu	Trp	Cys 150	Gln	Ala	Arg	Gly	Ile 155	Thr	Tyr	Thr	Arg	Tyr 160

65

CGC Arg 70	CGC Arg	TAC Tyr	ACC Thr	CCG Pro	GGC Gly 75	CGG Arg	AAG Lys	AAG Lys	TGG Trp	ATG Met 80	GAG Glu	GCC Ala	GCC Ala	GAG Glu	GCC Ala 85	533
CGG Arg	CGG Arg	CTG Leu	TTC Phe	TCC Ser 90	GCC Ala	ACG Thr	CTG Leu	CGC Arg	ACG Thr 95	CGG Arg	AAC Asn	CGG Arg	AAC Asn	CTG Leu 100	AGG Arg	581
GAC Asp	TTG Leu	CTG Leu	CCC Pro 105	GAC Asp	GAG Glu	GCA Ala	CAG Gln	CTG Leu 110	GCG Ala	CGC Arg	TAC Tyr	GGC Gly	CTG Leu 115	CCG Pro	GTC Val	629
TGG Trp	CGC Arg	ACG Thr 120	GAA Glu	GAG Glu	GAC Asp	GTG Val	GCA Ala 125	GCG Ala	GCC Ala	CTG Leu	GGC Gly	GTC Val 130	TCG Ser	GTG Val	GGC Gly	677
GTG Val 135	CTC Leu	CGC Arg	CAC His	TAC Tyr	AGC Ser	ATC Ile 140	CAC His	CGC Arg	CCG Pro	CGC Arg	GAG Glu 145	CGG Arg	GTG Val	CGG Arg	CAC His	725
TAC Tyr 150	GTG Val	ACC Thr	TTC Phe	GCC Ala	GTG Val 155	CCC Pro	AAG Lys	CGC Arg	TCC Ser	GGA Gly 160	GGC Gly	GTC Val	CGG Arg	CTG Leu	CTG Leu 165	773
CAT His	GCG Ala	CCC Pro	AAG Lys	CGG Arg 170	CGC Arg	CTG Leu	AAG Lys	GCC Ala	CTG Leu 175	CAA Gln	CGC Arg	CGG Arg	ATG Met	CTG Leu 180	GCG Ala	821
CTC Leu	CTG Leu	GTG Val	TCG Ser 185	AAG Lys	CTC Leu	CCC Pro	GTG Val	AGT Ser 190	CCA Pro	CAG Gln	GCC Ala	CAT His	GGC Gly 195	TTC Phe	GTG Val	869
CCC Pro	GGC Gly	CGC Arg 200	TCC Ser	ATC Ile	AAG Lys	ACG Thr	GGC Gly 205	GCC Ala	GCG Ala	CCG Pro	CAC His	GTG Val 210	GGC Gly	CGG Arg	CGG Arg	917
GTG Val 215	GTC Val	CTG Leu	AAG Lys	CTG Leu	GAC Asp	CTG Leu 220	AAG Lys	GAC Asp	TTC Phe	TTC Phe	CCC Pro 225	TCC Ser	GTC Val	ACC Thr	TTC Phe	965
GCG Ala 230	CGG Arg	GTG Val	CGA Arg	GGG Gly	CTG Leu 235	CTC Leu	ATC Ile	GCC Ala	CTG Leu	GGC Gly 240	TAC Tyr	GGC Gly	TAT Tyr	CCC Pro	GTG Val 245	1013
GCG Ala	GCC Ala	ACG Thr	CTC Leu	GCG Ala 250	GTG Val	CTG Leu	ATG Met	ACG Thr	GAG Glu 255	TCC Ser	GAG Glu	CGC Arg	CAG Gln	CCC Pro 260	GTG Val	1061
GAG Glu	CTG Leu	GAG Glu	GGC Gly 265	ATC Ile	CTC Leu	TTC Phe	CAC His	GTT Val 270	CCC Pro	GTG Val	GGC Gly	CCA Pro	CGC Arg	GTC Val 275	TGC Cys	1109

```
Ile Val Val
    210
```

CAC CTG CGA CAG GTC CGC CGG GAT GCG CGG CTG CTC CCC AAG GGC GTC 485  
His Leu Arg Gln Val Arg Arg Asp Ala Arg Leu Leu Pro Lys Gly Val



GTG CAG GGC GCC CCC ACG AGC CCC GCC CTG TGC AAC GCG GTG CTG CTG	1157
Val Gln Gly Ala Pro Thr Ser Pro Ala Leu Cys Asn Ala Val Leu Leu	
280 285 290	
CGA CTG GAC CGG CGG CTG GCG GGA CTG GCG CGT CGG TAC GGC TAC ACG	1205
Arg Leu Asp Arg Arg Leu Ala Gly Leu Ala Arg Arg Tyr Gly Tyr Thr	
295 300 305	
TAC ACG CGC TAC GCG GAT GAC CTC ACC TTC TCC GGC GAC GAC GTC ACG	1253
Tyr Thr Arg Tyr Ala Asp Asp Leu Thr Phe Ser Gly Asp Asp Val Thr	
310 315 320 325	
GCG CTG GAG CGA GTC CGC GCG CTG GCC GCG CGG TAC GTG CAG GAG GAA	1301
Ala Leu Glu Arg Val Arg Ala Leu Ala Ala Arg Tyr Val Gln Glu Glu	
330 335 340	
GGC TTC GAG GTC AAC CGC GAG AAG ACC CGC GTG CAG CGC CGG GGC GGT	1349
Gly Phe Glu Val Asn Arg Glu Lys Thr Arg Val Gln Arg Arg Gly Gly	
345 350 355	
GCC CAG CGC GTC ACT GGC GTC ACC GTG AAT ACG ACG CTG GGC TTG TCA	1397
Ala Gln Arg Val Thr Gly Val Thr Val Asn Thr Thr Leu Gly Leu Ser	
360 365 370	
CGC GAG GAG CGG CCG CGG CTC CGG GCG ATG CTG CAC CAG GAG GCG CGG	1445
Arg Glu Glu Arg Pro Arg Leu Arg Ala Met Leu His Gln Glu Ala Arg	
375 380 385	
TCG GAG GAC GTC GAG GCA CAC CGC GCG CAC CTC GAC GGC CTC CTG GCC	1493
Ser Glu Asp Val Glu Ala His Arg Ala His Leu Asp Gly Leu Leu Ala	
390 395 400 405	
TAC GTG AAG ATG CTC AAC CCG GAG CAG GCG GAG CGG CTC GCT CGC CGG	1541
Tyr Val Lys Met Leu Asn Pro Glu Gln Ala Glu Arg Leu Ala Arg Arg	
410 415 420	
CGC AAG CCG CGC GGG ACG TGAGCGAGGG CTCAGCTCCG GATGGGCCAG	1589
Arg Lys Pro Arg Gly Thr	
425	
GGCCTGTCAC GCGTCCCGGC CTCCAGTTG TCATGGCGGC CGTCCCAGTA C	1640

(2) INFORMATION FOR SEQ ID NO:38:

(i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 3060 base pairs
- (B) TYPE: nucleic acid
- (C) STRANDEDNESS: double
- (D) TOPOLOGY: linear

(ix) FEATURE:

- (A) NAME/KEY: CDS

(B) LOCATION: 763..2202

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:38:

CCCAC TTCCG GCGCTCGGGC TGC GCGAGGG CCCGTGCGAG CACATGATGG CGCTGCGGGCT 60  
CGTCCAGGTC CGGCACCGCG CCGAGCAGGA AGCACTGCGT CAGACCCCCG CGGGCCGCCA 120  
GCTCATCCGC GCGGAGACGC GCTCCTACGT GCGGCGCGAG CCCTCCGGCC AGGAGCAGGT 180  
GTACCGCGTC TCATTGGATG GGAAAGTGGT GGCGGTGGAG TGGGGCCCCC GCCAGGGGGA 240  
GTCCCGCCGG CAGAAGCTCT GGTTCGACAC GGACGCCGAG GCGCGCACCG CCTACTTCAC 300  
GCGCCTGGAG TCCTTGCCCG CGGAGGGATA TATCGATGCG GCTGCTTCAA TGATGTAGAA 360  
CACGCAAGCC ACGGGGCCGC GGGCGCGCGG CGGAAAGGCA GGTGCGACGG AACGACAGAC 420  
ACTCGTGCGA GCGACCGAGA GAGGTCCCAA GCCATCAGCC TCAGCGCCTC GAGCGCGAGA 480  
GCGGCGTTGC GCCGCTCTGG TTGAATTGCA GGACACTCTC CGCAAGGTAG CCTGTTCTTG 540  
GCTCTCTTCC CTCCGGTGAG TACCTCTCCG GCCGGGGAGC TGAACCAACG ACGCAACCGC 600  
CGTTTCCCCG GCCGGAGAGG TACTCACCGG AGGGGAGAGC CGGTGAGGCT ACCGTGCCCC 660  
AGGTGAGAAG GTGGTGCCCTT CGGGCCTCCC TCGACCGCTC GCGCTCCGTC GCCCTGCCCT 720  
GCCTCGCCCC CCCACCTTG CTCACCGGCG CCAGGAGCCG TC ATG ACC GCC AAG 774  
Met Thr Ala Lys  
1  
CTG GAG TCA CAC GTC CCC GCC GCG CCC CCC GTC TCC GCC GAG GCG CCC 822  
Leu Glu Ser His Val Pro Ala Ala Pro Pro Val Ser Ala Glu Ala Pro  
5 10 15 20  
GCC CCC ACC CGT CCC GAT GCC GCG AAG CAG GAG GCC CGC CGC GCC CAC 870  
Ala Pro Thr Arg Pro Asp Ala Ala Lys Gln Glu Ala Arg Arg Ala His  
25 30 35  
CAC GAG GCG CTG CGC CTG CGG TGG AAG GCC ATC GAA GAG GCG GGC GGC 918  
His Glu Ala Leu Arg Leu Arg Trp Lys Ala Ile Glu Glu Ala Gly Gly  
40 45 50  
ACG GAC GCC TGG GTG CGG CAG CAG CTG GTG GCC AAG GGC GTC GCG GCG 966  
Thr Asp Ala Trp Val Arg Gln Gln Leu Val Ala Lys Gly Val Ala Ala  
55 60 65  
GAA GAG GTG GAC TTC GAG TCG CTC AGC GAC AAG CAG AAG GCG GCC TGG 1014  
Glu Glu Val Asp Phe Glu Ser Leu Ser Asp Lys Gln Lys Ala Ala Trp  
70 75 80  
AAG GAG AAG AAG AAG GCC GAG GCC ACC GAG CGG CGC GCG CAG AAG CGC 1062

Lys 85	Glu	Lys	Lys	Lys	Ala 90	Glu	Ala	Thr	Glu	Arg 95	Arg	Ala	Gln	Lys	Arg 100	
CTG Leu	GCG Ala	TGG Trp	GAG Glu	GCC Ala	TGG Trp	AAG Lys	GCC Ala	ACG Thr	CAC His	ATC Ile	CAC His	CAC His	CTG Leu	GGC Gly	GTG Val	1110
GGG Gly	GTG Val	CAC His	TGG Trp	GAC Asp	GAG Glu	GCC Ala	GGA Gly	GGG Gly	CCG Pro	GAC Asp	AAG Lys	TTC Phe	GAC Asp	GTG Val	GCC Ala	1158
GGG Gly	CGC Arg	GAG Glu	GAG Glu	CGG Arg	GCC Ala	AAG Lys	GCC Ala	AAC Asn	GGC Gly	TTG Leu	CCG Pro	GAG Glu	GGG Gly	TTG Leu	GAC Asp	1206
TCG Ser	GTC Val	GAG Glu	GCG Ala	CTG Leu	GCC Ala	AAA Lys	GCG Ala	CTG Leu	GGC Gly	ATC Ile	TCC Ser	GTG Val	TCG Ser	CGC Arg	CTG Leu	1254
CGC Arg	TGG Trp	TTC Phe	TCC Ser	TTC Phe	CAC His	CGC Arg	GAG Glu	GTG Val	GAC Asp	ACG Thr	GGC Gly	ACG Thr	CAC His	TAC Tyr	CAG Gln	1302
ACG Thr	TGG Trp	GAG Glu	ATT Ile	CCG Pro	AAG Lys	CGG Arg	GAC Asp	GGC Gly	GGC Gly	AAG Lys	CGG Arg	ACG Thr	CTC Leu	ACC Thr	GCG Ala	1350
CCG Pro	AAG Lys	CGG Arg	GAG Glu	CTC Leu	AAG Lys	GCC Ala	GTG Val	CAG Gln	CGC Arg	TGG Trp	GTG Val	CTC Leu	GCG Ala	AAC Asn	GTG Val	1398
GTG Val	GAG Glu	CGG Arg	CTG Leu	CCG Pro	GTG Val	CAC His	GGG Gly	GCC Ala	GCG Ala	CAC His	GGC Gly	TTC Phe	GTG Val	GCG Ala	GGG Gly	1446
CGC Arg	TCC Ser	ATC Ile	CTC Leu	ACC Thr	AAC Asn	GCG Ala	CTG Leu	GCC Ala	CAC His	CAG Gln	GGC Gly	GCG Ala	GAC Asp	GTG Val	GTG Val	1494
GTG Val	AAG Lys	GTG Val	GAC Asp	ATG Met	AAG Lys	GAC Asp	TTC Phe	TTC Phe	CCT Pro	TCC Ser	GTG Val	ACG Thr	TGG Trp	CCC Pro	CGG Arg	1542
GTC Val	AAG Lys	GGA Gly	CTG Leu	CTG Leu	CGC Arg	AAG Lys	GGA Gly	GGA Gly	CTC Leu	CCG Pro	GAG Glu	AAC Asn	CTG Leu	GCG Ala	ACG Thr	1590
CTC Leu	CTG Leu	GCG Ala	CTG Leu	CTC Leu	TCC Ser	ACC Thr	GAG Glu	GCC Ala	CCG Pro	CGC Arg	GAG Glu	GTG Val	GTG Val	CGG Arg	TTC Phe	1638
CGG Arg	GGA Gly	GAG Glu	ACG Thr	CTG Leu	TAC Tyr	GTG Val	GCC Ala	AAG Lys	GGC Gly	CCT Pro	CGC Arg	GCG Ala	CTG Leu	CCC Pro	CAG Gln	1686

GGG	GCC	CCC	ACC	TCT	CCG	GCG	CTG	ACG	AAC	GCG	CTG	TGC	CTG	CGG	CTG	1734
Gly	Ala	Pro	Thr	Ser	Pro	Ala	Leu	Thr	Asn	Ala	Leu	Cys	Leu	Arg	Leu	
	310					315					320					
GAC	AAG	CGG	CTC	TCG	GCG	CTG	TCG	AAG	CGG	CTG	GGC	TTC	ACG	TAC	ACG	1782
Asp	Lys	Arg	Leu	Ser	Ala	Leu	Ser	Lys	Arg	Leu	Gly	Phe	Thr	Tyr	Thr	
325					330					335					340	
CGC	TAT	GCG	GAT	GAC	CTG	ACG	TTC	TCC	TGG	CGG	CGG	GCG	AAG	AAG	TCC	1830
Arg	Tyr	Ala	Asp	Asp	Leu	Thr	Phe	Ser	Trp	Arg	Arg	Ala	Lys	Lys	Ser	
				345					350					355		
CGG	CAG	AAG	GAA	CTC	CCC	CTG	GCG	GAT	GCG	CCG	GTG	GCG	CTG	CTC	CTG	1878
Arg	Gln	Lys	Glu	Leu	Pro	Leu	Ala	Asp	Ala	Pro	Val	Ala	Leu	Leu	Leu	
			360					365					370			
GCG	CGG	GTG	AAG	GGT	GTG	CTG	GAG	GCC	GAG	GGT	TTC	ACG	CTG	CAC	CCG	1926
Ala	Arg	Val	Lys	Gly	Val	Leu	Glu	Ala	Glu	Gly	Phe	Thr	Leu	His	Pro	
		375					380					385				
GAC	AAG	ACG	CGG	GTG	CAG	CGC	AAG	GGC	AGC	CGG	CAG	CGG	GTG	ACG	GGG	1974
Asp	Lys	Thr	Arg	Val	Gln	Arg	Lys	Gly	Ser	Arg	Gln	Arg	Val	Thr	Gly	
	390					395					400					
CTC	GTG	GTG	AAC	GAG	GCC	CCC	GAG	GGC	GTT	CCG	GGT	GCC	CGG	GTG	CCC	2022
Leu	Val	Val	Asn	Glu	Ala	Pro	Glu	Gly	Val	Pro	Gly	Ala	Arg	Val	Pro	
405					410					415					420	
CGC	GAT	GTG	GTG	CGG	CGG	CTG	CGC	GCG	GCG	ATC	CAC	AAC	CGG	GAG	CAG	2070
Arg	Asp	Val	Val	Arg	Arg	Leu	Arg	Ala	Ala	Ile	His	Asn	Arg	Glu	Gln	
				425				430						435		
GGC	AAG	CCC	GGC	CCC	ACC	GGG	GAG	ACG	CTG	GAG	CAG	CTC	AAG	GGG	CTC	2118
Gly	Lys	Pro	Gly	Pro	Thr	Gly	Glu	Thr	Leu	Glu	Gln	Leu	Lys	Gly	Leu	
			440					445					450			
GCG	GCC	TTC	CTT	CAC	ATG	ACG	GAC	GCG	GAG	AAG	GGC	CGC	GCC	TTC	CTG	2166
Ala	Ala	Phe	Leu	His	Met	Thr	Asp	Ala	Glu	Lys	Gly	Arg	Ala	Phe	Leu	
		455					460					465				
CGA	CGG	CTG	GAG	GCC	CTC	GAG	AAG	CGC	CAG	ACC	GCC	TGACCCTCAC				2212
Arg	Arg	Leu	Glu	Ala	Leu	Glu	Lys	Arg	Gln	Thr	Ala					
	470					475					480					
TGGTCGTCCG	GGGCATCGCA	GCGGGCGCCG	GGACGGACCG	TCACCCCCCA	GATCTCCATG											2272
CCATGCTGGG	GATTCTGGGC	GGTGAAGAAG	ACTTCCCAGC	CGAGACGGAC	GAAGCCCTGC											2332
GGATCCGATG	ACTCCTCGCC	CGGGGCGATC	TCCCGGAGGG	GCACCGTTCC	GACGTCCGTG											2392
CCATTGCTCA	CCCAGGGCTC	CCGGCCCCAG	CCTTGGGTGT	CCGCCGAGAA	GAAGAGCAGC											2452
CCGGAGATGG	CCGTCAGGTT	CTCCGGCGAC	GCATCCTCGG	GGCCCGGCGC	CAAATCCTTC											2512

250000-1-000000

AGCAGCAGGG TGCCCTTGGC GGTGCCATCG CTGGACCACA GCTCCCGGCC GTGGAGGCTG 2572  
TCACTCGCGG CGAAGTAGAG CATCCCATTG AGCGCCTTGA TGGCGCTGGG CGCCGAGCTG 2632  
TCCGGACCCG GCCAGATGTC CTTACCCCGG ACCGTGCCAT GCGACGTGCC ATCGCTGACC 2692  
CACAGCTCCT CGCCCTCGGG CTGGCCCCAG AACTCGGGCT CGCCTCCCCC GGCGCTGAAG 2752  
AAGATCTTCC CCCCAGAGCGC CGTGAGATCA TGC GGATAGA GGCCGGGGAA GAAGCGCAGC 2812  
TGCTCGGAGA CGGTGCCTCT GGAGCACCAC AGGCTGGCCT CGCCTTCGTC ATTGTGCGAGC 2872  
AGGAAGAAGA GCACCGAGTC CGCCGCGGTG AACGCGGAGA GGAAGTTGTC CTCGGGGCCC 2932  
GTGAAGACAG ACGTGGTGCT GGACAGCCCC AGGCTGCGCC AGATGAACAC CTCGTCATTG 2992  
ACGTTGGCCA CGAAGAAGAG CGCATCGCCG ACCCGGGTGA GCCGGCGCGG GCTGGAGCTG 3052  
CCGGGCAC 3060

(2) INFORMATION FOR SEQ ID NO:39:

(i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 2788 base pairs
- (B) TYPE: nucleic acid
- (C) STRANDEDNESS: double
- (D) TOPOLOGY: linear

(ix) FEATURE:

- (A) NAME/KEY: CDS
- (B) LOCATION: 2..103

(ix) FEATURE:

- (A) NAME/KEY: CDS
- (B) LOCATION: 707..1654

(ix) FEATURE:

- (A) NAME/KEY: CDS
- (B) LOCATION: 1644..2591

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:39:

T TTC GAG AAG CGC CAT ACC AAA CAG GGG ATA CAG ACC AAC CTG ACG 46  
Phe Glu Lys Arg His Thr Lys Gln Gly Ile Gln Thr Asn Leu Thr  
1 5 10 15  
CTG AAA GAG GAA AGC TAC GGC GAC TGG CTG CCG AAG TGC GAC GAC CCC 94  
Leu Lys Glu Glu Ser Tyr Gly Asp Trp Leu Pro Lys Cys Asp Asp Pro  
20 25 30  
GCA GCA ACA TAACCTCACT CAGACCGGCA ACAGCCGGTC TTTTCCTTTC 143

Ala Ala Thr

TGGCCATTGC CACAAGGTGA ACAATCCACT GTTCACCCCTT CACCGTTTAT TCACCCTTTA 203  
TCACTATGAA ATTATTAATA AAAAACCAGA GGTGAACAGT GTGAACAGTA AAACCTGAAA 263  
AAACTTTTTTA TCACCCCGCG CATCGCCCGA CTGGACAGAT CCAGAACGAG CAAAAATCAC 323  
AAAGGTGACG AGTCGACTGT TCACTCTTCA CCAACTCATC ACCACCTAAC CACATGATAT 383  
AAAATGATAA ATAATCGAGG TGAACAGTTA AATGCAAAAA AACTTTTTTCT CAGCTCTTGG 443  
ATAAAAGAAA ATTAATTCAC ATCAATAGCT TTCCTCTTGA ATCCTCTTGA GGTTTATGAG 503  
AGCGTAACAG AGCCAAACCT AGCATTTTAT GGGTTAATAG CCCATCGCGC ATGAGTCATG 563  
GTTTCGCCTA GTATTTTAGC TATGCCCCGTC GTTCAGTTTCG CTGAGCGGCG GCTGGGGGCC 623  
ACCGATCAGC GAACTGATCG ACGTGCTCAA GTAGGTTTGG CTCTTTTAGT CCTCTACCAT 683  
CAAGGTGCAT AAGGATATTC TCG ATG CTG ACT CAG CTA AAA AAA AAT GGT 733  
Met Leu Thr Gln Leu Lys Lys Asn Gly  
1 5  
ACT GAG GTA TCT AGA GCA ACC GCG TTA TTT TCA TCA TTC GTT GAA AAG 781  
Thr Glu Val Ser Arg Ala Thr Ala Leu Phe Ser Ser Phe Val Glu Lys  
10 15 20 25  
AAC AAA GTA AAA TGT CCT GGT AAT GTA AAA AAA TTC GTC TTT CTG TGT 829  
Asn Lys Val Lys Cys Pro Gly Asn Val Lys Lys Phe Val Phe Leu Cys  
30 35 40  
GGT GCT AAC AAA AAC AAT GGA GAA CCA TCA GCA AGA CGA TTG GAA TTA 877  
Gly Ala Asn Lys Asn Asn Gly Glu Pro Ser Ala Arg Arg Leu Glu Leu  
45 50 55  
ATA AAT TTT TCT GAA AGG TAT TTG AAT AAC TGT CAC TTT TTT CTT GCT 925  
Ile Asn Phe Ser Glu Arg Tyr Leu Asn Asn Cys His Phe Phe Leu Ala  
60 65 70  
GAA CTA GTT TTC AAA GAA TTA AGC ACC GAT GAA GAA TCA TTA TCT GAT 973  
Glu Leu Val Phe Lys Glu Leu Ser Thr Asp Glu Glu Ser Leu Ser Asp  
75 80 85  
AAT TTA TTA GAT ATC GAA GCT GAC TTA TCT AAA TTA GCT GAT CAT ATT 1021  
Asn Leu Leu Asp Ile Glu Ala Asp Leu Ser Lys Leu Ala Asp His Ile  
90 95 100 105  
ATC ATT GTT TTA GAA AGT TAT TCA TCT TTC ACG GAA CTT GGT GCA TTC 1069  
Ile Ile Val Leu Glu Ser Tyr Ser Ser Phe Thr Glu Leu Gly Ala Phe  
110 115 120  
GCA TAC AGC AAG CAA TTA CGC AAG AAA TTA ATA ATA GTT AAC AAT ACA 1117

Ala	Tyr	Ser	Lys 125	Gln	Leu	Arg	Lys	Lys 130	Leu	Ile	Ile	Val	Asn 135	Asn	Thr		
AAA	TTT	ATA	AAT	GAG	AAA	TCA	TTT	ATA	AAT	ATG	GGA	CCA	ATA	AAG	GCT	1165	
Lys	Phe	Ile	Asn 140	Glu	Lys	Ser	Phe 145	Ile	Asn	Met	Gly	Pro 150	Ile	Lys	Ala		
ATT	ACT	CAG	CAA	TCA	CAA	CAA	TCT	GGT	CAT	TTC	TTA	CAT	TAT	AAA	ATG	1213	
Ile	Thr	Gln	Gln	Ser	Gln	Gln	Ser 160	Gly	His	Phe	Leu 165	His	Tyr	Lys	Met		
ACA	GAA	GGT	ATT	GAA	AGT	ATA	GAG	CGC	TCT	GAT	GGG	ATT	GGC	GAA	ATA	1261	
Thr	Glu	Gly	Ile	Glu	Ser 175	Ile	Glu	Arg	Ser	Asp 180	Gly	Ile	Gly	Glu	Ile 185		
TTC	GAC	CCC	CTA	TAT	GAT	ATT	CTT	TCT	AAG	AAC	GAC	AGA	GCA	ATT	TCA	1309	
Phe	Asp	Pro	Leu 190	Tyr	Asp	Ile	Leu	Ser	Lys 195	Asn	Asp	Arg	Ala	Ile 200	Ser		
AGA	ACT	TTA	AAA	AAA	GAA	GAG	TTA	GAT	CCT	TCC	AGT	AAC	TTC	AAT	AAA	1357	
Arg	Thr	Leu	Lys 205	Lys	Glu	Glu	Leu 210	Asp	Pro	Ser	Ser	Asn 215	Phe	Asn	Lys		
GAC	TCA	GTA	CGA	TTT	ATT	CAT	GAC	GTA	ATT	TTT	GTA	TGT	GGT	CCT	TTG	1405	
Asp	Ser	Val 220	Arg	Phe	Ile	His	Asp 225	Val	Ile	Phe	Val	Cys 230	Gly	Pro	Leu		
CAA	CTT	AAT	GAA	CTC	ATC	GAA	ATA	ATC	ACA	AAA	ATA	TTT	GGC	ACA	GAA	1453	
Gln	Leu	Asn	Glu	Leu	Ile 240	Glu	Ile	Ile	Thr	Lys 245	Ile	Phe	Gly	Thr	Glu		
AGC	CAT	TAC	AAA	AAA	AAT	CTT	CTA	AAG	CAC	CTT	GGT	ATT	CTA	ATA	GCT	1501	
Ser	His	Tyr	Lys	Lys	Asn 255	Leu	Leu	Lys	His 260	Leu	Gly	Ile	Leu	Ile	Ala 265		
ATT	AGA	ATA	ATA	TCA	TGC	ACA	AAT	GGG	ATT	TAT	TAT	TCT	TTG	TAT	AAA	1549	
Ile	Arg	Ile	Ile 270	Ser	Cys	Thr	Asn	Gly 275	Ile	Tyr	Tyr	Ser	Leu	Tyr 280	Lys		
GAA	TAT	TAT	TTT	AAA	TAT	GAC	TTT	GAC	ATT	GAC	AAC	ATA	TCA	TCA	ATG	1597	
Glu	Tyr	Tyr	Phe 285	Lys	Tyr	Asp	Phe 290	Asp	Ile	Asp	Asn	Ile 295	Ser	Ser	Met		
TTT	AAA	GTT	TTT	TTC	CTC	AAG	AAC	AAG	CCA	GAA	AGG	ATG	AGG	GTA	TAT	1645	
Phe	Lys	Val 300	Phe	Phe	Leu	Lys	Asn 305	Lys	Pro	Glu	Arg	Met 310	Arg	Val	Tyr		
GAG	AAT	ATA	TAGCCTAATT GATTCTCAGA CATTGATGAC TAAGGGATT T													1694	
Glu	Asn	Ile 315															
GCTTCTGAAG	TAATGCGATC ACCTGAGCCG CCAAAAAAAT GGGATATAGC TAAGAAAAAA															1754	
GGAGGTATGA	GAACAATT TAACCCGTCA TCAAAAGTTA AATTAATTCA ATATTGGTTA															1814	

ATGAATAATG TTTTTCGAA GCTCCCAATG CATAATGCTG CATATGCATT TGTAAAAAC 1874  
CGATCAATAA AAAGCAATGC TTTATTACAT GCCGAATCAA AGAATAAGTA TTATGTGAAA 1934  
ATAGATCTCA AAGATTTTTT CCCTTCAATA AAATTTACTG ATTTTGAGTA CGCATTCACT 1994  
CGTTATCGAG ATCGCATTGA ATTTACTACA GAATATGATA AGGAGTTACT ACAACTTATA 2054  
AAAACGATCT GCTTTATATC AGATAGCACT CTCCCTATCG GGTTTCCTAC ATCTCCATTA 2114  
ATTGCAAACT TTGTGGCAAG AGAACTTGAT GAAAAACTGA CGCAAAACT AAATGCAATT 2174  
GATAAACTTA ATGCCACTTA TACACGATAT GCTGATGATA TTATTGTCTC TACAAATATG 2234  
AAAGGGGCTA GCAAATTAAT TCTGGATTGT TTTAAAAGAA CAATGAAAGA GATTGGTCCA 2294  
GACTTTAAAA TTAACATTAA AAAATTTAAG ATTTGTAGTG CTTCGGGAGG AAGTATAGTA 2354  
GTTACCGGAT TGAAAGTTTG CCACGATTTT CATATTACAT TACATAGATC AATGAAAGAT 2414  
AAAATAAGAT TGCATCTTTC TCTTTTATCA AAGGGCATAT TAAAAGATGA AGATCATAAT 2474  
AACTTTCTG GTTATATTGC TTATGCAAAA GATATAGACC CTCATTTTTA TACAAAACCTG 2534  
AACAGAAAAT ATTTTCAAGA AATAAAATGG ATTCAGAATC TCCACAACAA AGTTGAATAA 2594  
ACTTTATATT TTGGATGCAC CCCAATAACT TCATTGATTA AATTGGGAAC AATATAGGCT 2654  
TTTCAGGATG ACCTACACTC TAGAGAATGT GTATACAAAA GTGTATAAGT TATTTTCAAA 2714  
CCTATATAAA ATACAGCAAA ATCAATGCAT TGGCGGCATT TTACCACTCC TGTGATCTTC 2774  
CGCCAAAATG CCTC 2788

(2) INFORMATION FOR SEQ ID NO:40:

(i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 316 amino acids
- (B) TYPE: amino acid
- (D) TOPOLOGY: linear

(ii) MOLECULE TYPE: protein

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:40:

Met	Arg	Ile	Tyr	Ser	Leu	Ile	Asp	Ser	Gln	Thr	Leu	Met	Thr	Lys	Gly
1				5					10					15	
Phe	Ala	Ser	Glu	Val	Met	Arg	Ser	Pro	Glu	Pro	Pro	Lys	Lys	Trp	Asp
			20					25						30	
Ile	Ala	Lys	Lys	Lys	Gly	Gly	Met	Arg	Thr	Ile	Tyr	His	Pro	Ser	Ser
		35					40						45		



Lys Val Lys Leu Ile Gln Tyr Trp Leu Met Asn Asn Val Phe Ser Lys  
 50 55 60  
 Leu Pro Met His Asn Ala Ala Tyr Ala Phe Val Lys Asn Arg Ser Ile  
 65 70 75 80  
 Lys Ser Asn Ala Leu Leu His Ala Glu Ser Lys Asn Lys Tyr Tyr Val  
 85 90 95  
 Lys Ile Asp Leu Lys Asp Phe Phe Pro Ser Ile Lys Phe Thr Asp Phe  
 100 105 110  
 Glu Tyr Ala Phe Thr Arg Tyr Arg Asp Arg Ile Glu Phe Thr Thr Glu  
 115 120 125  
 Tyr Asp Lys Glu Leu Leu Gln Leu Ile Lys Thr Ile Cys Phe Ile Ser  
 130 135 140  
 Asp Ser Thr Leu Pro Ile Gly Phe Pro Thr Ser Pro Leu Ile Ala Asn  
 145 150 155 160  
 Phe Val Ala Arg Glu Leu Asp Glu Lys Leu Thr Gln Lys Leu Asn Ala  
 165 170 175  
 Ile Asp Lys Leu Asn Ala Thr Tyr Thr Arg Tyr Ala Asp Asp Ile Ile  
 180 185 190  
 Val Ser Thr Asn Met Lys Gly Ala Ser Lys Leu Ile Leu Asp Cys Phe  
 195 200 205  
 Lys Arg Thr Met Lys Glu Ile Gly Pro Asp Phe Lys Ile Asn Ile Lys  
 210 215 220  
 Lys Phe Lys Ile Cys Ser Ala Ser Gly Gly Ser Ile Val Val Thr Gly  
 225 230 235 240  
 Leu Lys Val Cys His Asp Phe His Ile Thr Leu His Arg Ser Met Lys  
 245 250 255  
 Asp Lys Ile Arg Leu His Leu Ser Leu Leu Ser Lys Gly Ile Leu Lys  
 260 265 270  
 Asp Glu Asp His Asn Lys Leu Ser Gly Tyr Ile Ala Tyr Ala Lys Asp  
 275 280 285  
 Ile Asp Pro His Phe Tyr Thr Lys Leu Asn Arg Lys Tyr Phe Gln Glu  
 290 295 300  
 Ile Lys Trp Ile Gln Asn Leu His Asn Lys Val Glu  
 305 310 315

(2) INFORMATION FOR SEQ ID NO:41:

(i) SEQUENCE CHARACTERISTICS:

(A) LENGTH: 1602 base pairs  
 (B) TYPE: nucleic acid  
 (C) STRANDEDNESS: double  
 (D) TOPOLOGY: linear

(ix) FEATURE:

(A) NAME/KEY: CDS  
 (B) LOCATION: 548..1507

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:41:

TGGCATCTAT TAAGAAGGTT AGGAAAGAAA ATAAAGTATC AAAAGATATT GGAAATATAT 60  
 TATACGCAGA GCGTTTCTAT TGCCTTGTAT CTATTTACTG GATAGTGTCA ACTACCGCAC 120  
 ACTGTGTGAA CTAGCTTTTA AAGCGATAAA GCAAGATGAT GTTTTATCTA AAATTATTGT 180  
 TAGATCCGTT GTTTCTCGTC TAATAAATGA ACGAAAAATA CTTCAAATGA CTGATGGTTA 240  
 TCAGGTCACT GCTTTGGGGG CTAGCTATGT TAGGAGCGTC TTTGATAGAA AGACACTTGA 300  
 CCGATTGCGG CTTGAGATTA TGAATTTTGA AAACCGTAGA AAATCAACAT TTAACATGA 360  
 TAAGATTCCG TATGCGCACC CTTAGCGAGA GGTTTATCAT TAAGGTCAAC CTCTGGATGT 420  
 TGTTTCGGCA TCCTGCATTG AATCTGAGTT ACTGTCTGTT TTCCTTGTTG GAACGGAGAG 480  
 CATCGCCTGA TGCTCTCCGA GCCAACCAGG AAACCCGTTT TTTCTGACGT AAGGGTGCGC 540  
 AACTTTC ATG AAA TCC GCT GAA TAT TTG AAC ACT TTT AGA TTG AGA AAT 589  
 Met Lys Ser Ala Glu Tyr Leu Asn Thr Phe Arg Leu Arg Asn  
 1 5 10  
 CTC GGC CTA CCT GTC ATG AAC AAT TTG CAT GAC ATG TCT AAG GCG ACT 637  
 Leu Gly Leu Pro Val Met Asn Asn Leu His Asp Met Ser Lys Ala Thr  
 15 20 25 30  
 CGC ATA TCT GTT GAA ACA CTT CGG TTG TTA ATC TAT ACA GCT GAT TTT 685  
 Arg Ile Ser Val Glu Thr Leu Arg Leu Leu Ile Tyr Thr Ala Asp Phe  
 35 40 45  
 CGC TAT AGG ATC TAC ACT GTA GAA AAG AAA GGC CCA GAG AAG AGA ATG 733  
 Arg Tyr Arg Ile Tyr Thr Val Glu Lys Lys Gly Pro Glu Lys Arg Met  
 50 55 60  
 AGA ACC ATT TAC CAA CCT TCT CGA GAA CTT AAA GCC TTA CAA GGA TGG 781  
 Arg Thr Ile Tyr Gln Pro Ser Arg Glu Leu Lys Ala Leu Gln Gly Trp  
 65 70 75  
 GTT CTA CGT AAC ATT TTA GAT AAA CTG TCG TCA TCT CCT TTT TCT ATT 829  
 Val Leu Arg Asn Ile Leu Asp Lys Leu Ser Ser Ser Pro Phe Ser Ile  
 80 85 90

GGA Gly 95	TTT Phe	GAA Glu	AAG Lys	CAC His	CAA Gln 100	TCT Ser	ATT Ile	TTG Leu	AAT Asn	AAT Asn 105	GCT Ala	ACC Thr	CCG Pro	CAT His	ATT Ile 110	877
GGG Gly	GCA Ala	AAC Asn	TTT Phe	ATA Ile 115	CTG Leu	AAT Asn	ATT Ile	GAT Asp	TTG Leu 120	GAG Glu	GAT Asp	TTT Phe	TTC Phe	CCA Pro 125	AGT Ser	925
TTA Leu	ACT Thr	GCT Ala	AAC Asn 130	AAA Lys	GTT Val	TTT Phe	GGA Gly	GTG Val 135	TTC Phe	CAT His	TCT Ser	CTT Leu	GGT Gly 140	TAT Tyr	AAT Asn	973
CGA Arg	CTA Leu	ATA Ile 145	TCT Ser	TCA Ser	GTT Val	TTG Leu	ACA Thr 150	AAA Lys	ATA Ile	TGT Cys	TGT Cys	TAT Tyr 155	AAA Lys	AAT Asn	CTG Leu	1021
CTA Leu 160	CCA Pro	CAA Gln	GGT Gly	GCT Ala	CCA Pro	TCA Ser 165	TCA Ser	CCT Pro	AAA Lys	TTA Leu	GCT Ala 170	AAT Asn	CTA Leu	ATA Ile	TGT Cys	1069
TCT Ser 175	AAA Lys	CTT Leu	GAT Asp	TAT Tyr	CGT Arg 180	ATT Ile	CAG Gln	GGT Gly	TAT Tyr	GCA Ala 185	GGT Gly	AGT Ser	CGG Arg	GGC Gly	TTG Leu 190	1117
ATA Ile	TAT Tyr	ACG Thr	AGA Arg	TAT Tyr 195	GCC Ala	GAT Asp	GAT Asp	CTC Leu	ACC Thr 200	TTA Leu	TCT Ser	GCA Ala	CAG Gln	TCT Ser 205	ATG Met	1165
AAA Lys	AAG Lys	GTT Val 210	GTT Val	AAA Lys	GCA Ala	CGT Arg	GAT Asp	TTT Phe 215	TTA Leu	TTT Phe	TCT Ser	ATA Ile	ATC Ile 220	CCA Pro	AGT Ser	1213
GAA Glu	GGA Gly	TTG Leu 225	GTT Val	ATT Ile	AAC Asn	TCA Ser	AAA Lys 230	AAA Lys	ACT Thr	TGT Cys	ATT Ile	AGT Ser 235	GGG Gly	CCT Pro	CGT Arg	1261
AGT Ser 240	CAG Gln	AGG Arg	AAA Lys	GTT Val	ACA Thr	GGT Gly 245	TTA Leu	GTT Val	ATT Ile	TCA Ser	CAA Gln 250	GAG Glu	AAA Lys	GTT Val	GGG Gly	1309
ATA Ile 255	GGT Gly	AGA Arg	GAA Glu	AAA Lys 260	TAT Tyr	AAA Lys	GAA Glu	ATT Ile	AGA Arg	GCA Ala 265	AAG Lys	ATA Ile	CAT His	CAT His	ATA Ile 270	1357
TTT Phe	TGC Cys	GGT Gly	AAG Lys	TCT Ser 275	TCT Ser	GAG Glu	ATA Ile	GAA Glu	CAC His 280	GTT Val	AGG Arg	GGA Gly	TGG Trp	TTG Leu 285	TCA Ser	1405
TTT Phe	ATT Ile	TTA Leu	AGT Ser 290	GTG Val	GAT Asp	TCA Ser	AAA Lys	AGC Ser 295	CAT His	AGG Arg	AGA Arg	TTA Leu	ATA Ile 300	ACT Thr	TAT Tyr	1453
ATT Ile	AGC Ser	AAA Lys	TTA Leu	GAA Glu	AAA Lys	AAA Lys	TAT Tyr	GGA Gly	AAG Lys	AAC Asn	CCT Pro	TTA Leu	AAT Asn	AAA Lys	GCG Ala	1501

305

310

315

AAG ACC TAATGGTCTT CGTTTTAAAA CTAAAGCTCA TAGGTTGAAA AATTGAGCAC 1557  
Lys Thr  
320

TTCTTCGTCC AACCAAGTTAT TTAGTTCCTG CAATCGTTTC TGCAG 1602

(2) INFORMATION FOR SEO ID NO:42:

(i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 1540 base pairs  
(B) TYPE: nucleic acid  
(C) STRANDEDNESS: double  
(D) TOPOLOGY: linear

(ix) FEATURE:

- (A) NAME/KEY: CDS  
(B) LOCATION: 396..1352

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:42:

TCACCCTGAA AGACCTGATT GCTTACCTGG AAGAGAAGCC GGAAATGGCG GAACATCTGG 60

CGGCGGTTAA GGCCTATCGC GAAGAGTTCG GCGTTTAAAA ATATGCGCTG TGCAGGGTTT 120

TTGCTGTGCG CAGCGTGATG CGCTTCAAGA TATCGTGTTA ATCTGCTTTC GCCAGCAGTG 180

GCAATAGCGT TTCCGGCCTT TTGTGCCGGG AGGGTCGGCG AGTCGCTGAC TTAACGCCAG 240

TAGTATGTCC ATATACCCAA AGTCGCTTCA TTGTACCTGA GTACGCTTCG CGTACGTCGC 300

GCTGACGCGC TCAGTACAGT TACGCGCCTT CGGGATGGTT TAATGGTATT GCCGCTGTTG 360

GCGCCTCTTT TGGCCGCCGT GATGTGGAGA GTGGA ATG GAT GCT ACC CGG ACA 413  
Met Asp Ala Thr Arg Thr  
1 5

ACC CTT CTG GCG CTC GAT TTG TTC GGC TCG CCG GGC TGG AGC GCC GAT 461  
Thr Leu Leu Ala Leu Asp Leu Phe Gly Ser Pro Gly Trp Ser Ala Asp  
10 15 20

AAA GAA ATA CAG CGA CTG CAT GCG CTC AGT AAT CAT GCC GGA CGC CAT 509  
Lys Glu Ile Gln Arg Leu His Ala Leu Ser Asn His Ala Gly Arg His  
25 30 35

TAC CGA CGC ATT ATT CTT TCT AAA CGC CAC GGT GGT CAG CGG CTG GTG 557  
Tyr Arg Arg Ile Ile Leu Ser Lys Arg His Gly Gly Gln Arg Leu Val  
40 45 50

TTA GCC CCT GAT TAC TTG CTC AAA ACC GTA CAG CGC AAC ATT CTT AAG 605

Leu 55	Ala	Pro	Asp	Tyr	Leu 60	Leu	Lys	Thr	Val	Gln 65	Arg	Asn	Ile	Leu	Lys 70	
AAC Asn	GTC Val	CTT Leu	TCA Ser	CAA Gln 75	TTT Phe	CCG Pro	CTT Leu	TCC Ser	CCT Pro 80	TTT Phe	GCT Ala	ACA Thr	GCC Ala	TAC Tyr 85	CGA Arg	653
CCA Pro	GGT Gly	TGC Cys 90	CCA Pro	ATC Ile	GTC Val	AGC Ser	AAC Asn	GCG Ala 95	CAG Gln	CCA Pro	CAC His	TGC Cys 100	CAA Gln 100	CAG Gln	CCG Pro	701
CAG Gln	ATC Ile	CTG Leu 105	AAA Lys	CTC Leu	GAT Asp	ATC Ile	GAA Glu 110	AAC Asn	TTT Phe	TTC Phe	GAT Asp	AGC Ser 115	ATT Ile	AGC Ser	TGG Trp	749
TTA Leu 120	CAG Gln	GTC Val	TGG Trp	CGT Arg	GTG Val	TTT Phe 125	CGC Arg	CAG Gln	GCC Ala	CAG Gln 130	TTG Leu	CCA Pro	CGT Arg	AAT Asn	GTG Val	797
GTA Val 135	ACC Thr	ATG Met	CTG Leu	ACC Thr	TGG Trp 140	ATT Ile	TGT Cys	TGT Cys	TAT Tyr	AAC Asn 145	GAC Asp	GCG Ala	TTA Leu	CCG Pro	CAG Gln 150	845
GGG Gly	GCA Ala	CCA Pro	ACT Thr	TCG Ser 155	CCA Pro	GCC Ala	ATT Ile	TCC Ser	AAT Asn 160	CTT Leu	GTG Val	ATG Met	CGC Arg 165	CGT Arg	TTT Phe	893
GAT Asp	GAA Glu	CGC Arg	ATA Ile 170	GGG Gly	GAA Glu	TGG Trp	TGT Cys	CAG Gln 175	GCT Ala	CGG Arg	GGA Gly	ATT Ile 180	ACC Thr 180	TAC Tyr	ACC Thr	941
CGC Arg	TAC Tyr	TGC Cys 185	GAT Asp	GAC Asp	ATG Met	ACC Thr	TTT Phe 190	TCA Ser	GGT Gly	CAC His	TTC Phe 195	AAT Asn 195	GCC Ala	CGC Arg	CAG Gln	989
GTT Val 200	AAA Lys	AAT Asn	AAA Lys	GTG Val	TGC Cys	GGA Gly 205	TTG Leu	TTA Leu	GCG Ala	GAG Glu 210	CTG Leu	GGC Gly	CTG Leu	AGC Ser	CTC Leu	1037
AAT Asn 215	AAA Lys	CGC Arg	AAA Lys	GGC Gly	TGC Cys 220	CTG Leu	ATA Ile	GCT Ala	GCC Ala	TGT Cys 225	AAG Lys	CGC Arg	CAG Gln	CAA Gln	GTA Val 230	1085
ACC Thr	GGG Gly	ATT Ile	GTT Val 235	GTT Val	AAT Asn	CAC His	AAG Lys	CCA Pro	CAG Gln 240	CTT Leu	GCC Ala	CGT Arg	GAA Glu 245	GCG Ala	CGC Arg	1133
CGG Arg	GCG Ala	CTG Leu	CGT Arg 250	CAG Gln	GAG Glu	GTG Val	CAT His 255	TTG Leu	TGC Cys	CAA Gln	AAA Lys	TAT Tyr 260	GGC Gly 260	GTT Val	ATT Ile	1181
TCG Ser	CAT His	CTT Leu 265	AGT Ser	CAT His	CGT Arg	GGT Gly 270	GAA Glu 270	CTT Leu	GAT Asp	CCT Pro	TCT Ser	GGC Gly 275	GAT Asp	CTC Leu	CAC His	1229

GCA CAG GCA ACG GCG TAT CTT TAT GCT TTG CAG GGA AGA ATA AAC TGG	1277
Ala Gln Ala Thr Ala Tyr Leu Tyr Ala Leu Gln Gly Arg Ile Asn Trp	
280 285 290	
TTA TTG CAA ATC AAC CCT GAG GAT GAG GCC TTT CAA CAG GCG AGA GAG	1325
Leu Leu Gln Ile Asn Pro Glu Asp Glu Ala Phe Gln Gln Ala Arg Glu	
295 300 305 310	
AGT GTA AAG CGA ATG CTG GTT GCA TGG TAAGAAAAGC GTCAGGCAGA	1372
Ser Val Lys Arg Met Leu Val Ala Trp	
315	
CGTTTCTGCC TGACCGTTTA GGGGAGAATT ACTGCAACTG CGCGGCAATT AGCGGCCAGC	1432
GGGCGTCAAA ATCATCCGTC GGGCGGTATT TAAACTCGCT GCGGACAAAA CGTGACAGCA	1492
TACCTTCACA GAAGGCCAGG ATCTGGCTTG CCAGCAGGGT TTCATCGG	1540

(2) INFORMATION FOR SEQ ID NO:43:

- (i) SEQUENCE CHARACTERISTICS:
  - (A) LENGTH: 4 amino acids
  - (B) TYPE: amino acid
  - (D) TOPOLOGY: linear

(ii) MOLECULE TYPE: protein

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:43:

Tyr Xaa Asp Asp  
1 4

(2) INFORMATION FOR SEQ ID NO:44:

- (i) SEQUENCE CHARACTERISTICS:
  - (A) LENGTH: 4 amino acids
  - (B) TYPE: amino acid
  - (D) TOPOLOGY: linear

(ii) MOLECULE TYPE: protein

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:44:

Ser Xaa Xaa Xaa  
1 4

(2) INFORMATION FOR SEQ ID NO:45:

(i) SEQUENCE CHARACTERISTICS:

(A) LENGTH: 4 amino acids

(B) TYPE: amino acid

(D) ~~TOPOLOGY~~: linear

(ii) MOLECULE TYPE: protein

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:45:

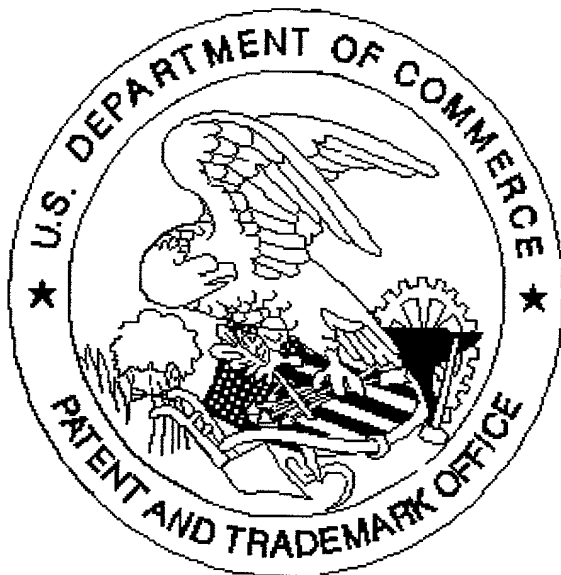
~~Xaa~~ Val Thr Gly  
1 4

1

4

**SECRET**

United States Patent & Trademark Office  
Office of Initial Patent Examination -- Scanning Division



Application deficiencies found during scanning:

☐ Page(s) \_\_\_\_\_ of \_\_\_\_\_ were not present  
for scanning. (Document title)

☐ Page(s) \_\_\_\_\_ of \_\_\_\_\_ were not present  
for scanning. (Document title)

☒ *Scanned copy is best available.* Drawings are dark, and  
there are lines in specification  
and sequence.